

LB

1131

.H655

Res. Lib.

The University of Chicago
Libraries

RELEASED
IBU:GAT



GIFT OF

Karl John Holzinger

THE UNIVERSITY OF CHICAGO

INDEXING A MENTAL CHARACTERISTIC

A DISSERTATION
SUBMITTED TO THE FACULTY
OF THE GRADUATE SCHOOL OF ARTS AND LITERATURE
IN CANDIDACY FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

DEPARTMENT OF EDUCATION

BY
KARL JOHN HOLZINGER

CHICAGO, ILLINOIS
SEPTEMBER, 1922

7175

RECEIVED

164:CAT

164:CAT

LB1131
H655



TABLE OF CONTENTS

Section 1. Introduction

- a. Mental and Physical Measurement
- b. A General Statistical Theory
- c. A Supposititious Example
- d. Indexing a Mental Characteristic

Section 2. The Groups and Data Studied

- a. The Groups Studied
- b. The Tests Used

Section 3. Test Administration and Scoring

- a. Primary and Secondary Index Variables
- b. The Difficulty Factor Eliminated
- c. Methods of Administering the Tests
- d. Authors' Plans of Scoring

Section 4. Methods Employed

Part I. ANALYSIS OF THE INDEX VARIABLES BY WHOLE SCALES

A. Comparative Validity and Reliability of the Indexes

Section 5. Relationships between Index Variables on the Same Scales

Section 6. Relationships between Scales by the Same Index Variable

Section 7. The Reliability of a Scale by Different Index Variables

Section 8. Correlations with Age

Section 9. Correlations with School Marks

Section 10. Summary of Reliability of Indexes for Whole Scales

B. The Discriminative Capacity of the Indexes

Section 11. Capacity of the Indexes to Discriminate between
Individuals

Section 12. Capacity of the Indexes to Discriminate between
Groups

Section 13. Practice Effect with Repetition

Part II. ANALYSIS OF THE INDEX VARIABLES BY COMPONENT TESTS

Section 14. Intercorrelations of Variables for the Otis
Components

Section 15. Correlations of the Otis Components with Age

Section 16. The Application of Reliability Formulae to the
Component Tests

Section 17. Summary of Analysis by Components

Part III. SCORING FORMULAE

Section 18. The Linear Form, $S = aR + cW$

a. Formulae with Highest Validity

b. Limitations in the Use of the Formula, $S = aR + cW$

c. Use of the Formula, $S = R + CW$

Section 19. Simple Ratios

a. The Correlation between Speed and Accuracy

b. The Validity of Simple Ratios as Scoring Indexes

c. The Reliability of Ratios as Scoring Indexes

Section 20. General Summary

Appendix A. Correlation Tables for Reliability Coefficients

Appendix B. Theorems Relating to Correlations

Bibliography

Section 1. Introduction

A. Mental and Physical Measurement

Most of the difficulties in mental measurement are due to the fact that the method is necessarily indirect. This indirectness comes not only in the traits measured, but also in the precision of the result. Any physical measurements on the other hand are direct, usually possible a clear-or definition of the traits measured and greater accuracy of their determination. The case of a boy to be weighed and also tested for intelligence will bring out this contrast.

If the boy is placed on the physical scale and weighed, the measurement is recorded directly in objectively defined units, accuracy of determination depending upon the instrument and the individual making the measurement. The reaction of the boy himself at the time of the weighing will have no effect on the result. In case a mental measurement of intelligence is required, it is necessary to submit to the boy a series of questions to which he must respond before a score for the mental trait may be obtained. This very indirectness of approach leads to ambiguity regarding the trait measured, and to uncertainty in the precision of the score. Thus in passing from physical to mental measurement the role of the boy has shifted from a passive to an active one, thereby complicating the entire procedure.

It may be stated without doubt, that this problem in the case of mental measurements has very closely paralleled

the physical method. Units have been defined and scales constructed. These units for the most part are functions of group variability on particular types of material, and consequently, less in accuracy and differences when applied in individual measurement. The parallelism in method has even been carried so far as to attempt to determine "zero points" for certain mental abilities so that "just not any" amounts of the traits could be used as reference points. In the case of the measurement of heat, the zero on the Fahrenheit thermometer does not mean "just no heat," but nevertheless serves as a convenient reference point. Similarly a good many of the former psychology zero points in mental measurement have disappeared, or have moved up to the position where they belong.

b. A General Statistical Theory

The difficulties briefly sketched in the foregoing paragraphs support the desirability of some more general method of approach to the problem of mental measurement, and indeed such a method has been implicitly used in some of the later work on test reliability and validity. In order to fix ideas, therefore, some in turn will now be defined so that their meaning will be clear throughout the subsequent discussion.

The term characteristic will be used to denote the physical, mental, or social traits of individuals of a group have in common. The group may consist of a number of persons or things each of which must possess the characteristic in question before any statistical study is possible. A knowledge

certain height, intelligence, and wealth furnishes an example of an individual with the three types of characteristics commonly studied.

The phase of a characteristic may be briefly described as the status of an individual with respect to the characteristic. This conception is an important one because it is introduced for the purpose of distinguishing a particular thing from the number that may be attached to it. Phases may be numerically or verbally expressed. Thus in describing the characteristics height and political affiliation of a certain man, the number 68 may be attached to the phase in height and the word "Republican" to the phase in political affiliation. It is conceivable that a numerical scheme for the latter might be worked out, but the phases of political affiliation per se would of course remain unchanged.

In order that a characteristic be numerically indexed, it is desirable that its phases be arranged in some order e.g. like the points on a line. If the linear arrangement be made, the trait may be termed a linear ordinal characteristic. For characteristics such as height, an infinite number of phases are assumed, indexed by the real number system (dense set). In the case of such characteristics as size of school class, however, the number of phases is finite, and the indexing is accomplished by assigning only integers. The distinction is essentially that between continuous and discrete series, the continuity and discreteness appearing in phase.

Finally an index variable will be described as a quan-

... ..
... ..
... ..

tity whose values are in one-to-one correspondence with the phases of the characteristics indexed. Before proceeding farther, the meaning of these various terms will be illustrated by means of an artificial example.

c. A Supposititious Example

Consider a set of 90 cubes of homogeneous material. The problem is to describe these cubes by ordinary statistical procedure. Assuming that the size of the cube is the characteristic to be studied, some mode of indexing or index variable must be adopted. Taking the edge as a first choice, the distribution may be given as follows:

edge	frequency
1	10
2	20
3	30
4	20
5	10

The mean edge is clearly 3, with corresponding face area 9 and volume 27 i.e. $\bar{x} = \bar{e}$; $\bar{y} = \bar{e}^2$.

A second mode of indexing by the area of a face gives the distribution:

area	frequency	fa
1	10	10
4	20	80
9	30	270
16	20	320
25	10	250
		<hr/> 730

The mean face area is now $17\frac{1}{3}$ with a corresponding edge and volume approximately 4.2 and 74.2 respectively.

Again indexing by volume gives:

volume	frequency	fv
1	10	10
8	30	180
27	30	310
64	30	1200
125	10	1250
		<u>5510</u>

The mean volume is 32, with the corresponding edge and area approximately 3.4 and 11.5. These results may be set forth in summary form as follows:

TABLE 1.- MEANS AND CORRESPONDING VALUES FOR CUBES INDEXED BY
EDGE, AREA, AND VOLUME

Mode of Indexing or Index Variable	Mean and Corresponding Indexes		
	Edge	Area	Volume
Edge	3	9	27
Area	3.2	<u>11.5</u>	32.2
Volume	3.4	11.5	<u>39</u>

Inspection of these figures reveals a complete inconsistency under the three modes of indexing. Moreover the three distributions are quite different. The frequency polygon according to edge is symmetrical, while the distributions for area and volume are skewed toward the smaller values of the variable. It is also clear that the cubes remained at the same phases in the characteristic size but that various modes of indexing gave inconsistent results. The above example then illustrates the fact that although a set of things be unaltered in phase, the form of the distribution and the statistical constants depend upon the particular index variable employed.

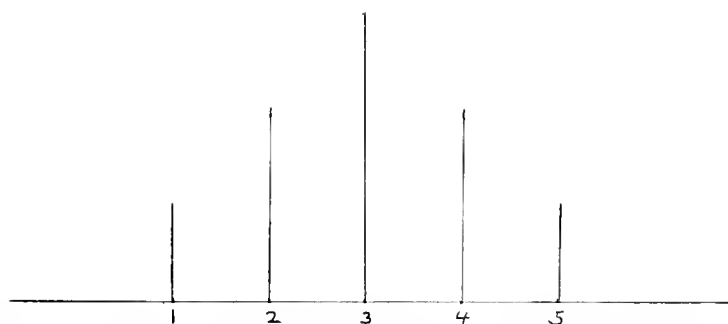
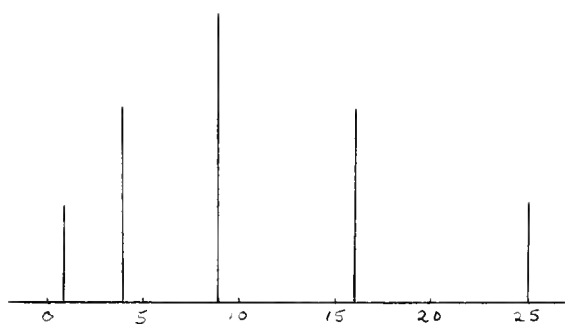
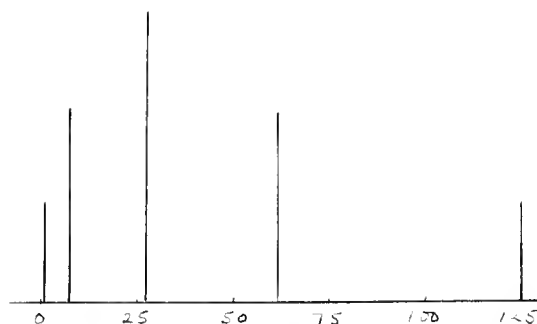


Figure 1.- Distribution of Cubes by Edge

Figure 2.- Distribution of
Cubes by AreaFigure 3.- Distribution of
Cubes by Volume

It may be well to point out briefly that the inconsistencies in the above example are due to the relationships between the index variables; i. e.

$$V = e^3$$

$$A = e^2$$

It will be sufficient to note here that similar inconsistencies will arise whenever the relationship between the index variables is other than linear ($y = ax + b$). In the latter case the statistical constants will be merely affected by propor-

tionality factors. Thus when height is measured in inches and again in centimeters, the distributions will be similar and all constants easily converted by the simple linear relationship,

$$1 \text{ inch} = 2.54 \text{ centimeters}$$

d. Indexing a Mental Characteristic

In mental measurement the method is to set a body of material before a child and elicit certain responses from him. These responses are then recorded and combined in various ways to produce what is known as a score. Moreover these responses exhibit considerable variety under different modes of administering the tests. Two boys took the Terman Group Test of Mental Ability which is administered under the plan of keeping the time constant. Their responses may be set forth briefly as follows:

	Score (Author)	Right	Attempted	Wrong	Accuracy
John	130	128	175	50	.71
Henry	150	140	150	10	.93

The problem of indexing here is ^{not} unlike that for the case of the cubes. Intelligence is the characteristic to be indexed, and this is possible by using the Rights, Wrongs, Attempts, or some combination of these index variables in the form of a score. It will be shown later that for the particular scale in question the author's score is a rather complicated function of the variables Right and Wrong, but for the present it is sufficient to note the possibility of such a characteristic

being indexed in several ways.

Questions immediately arise as to the best mode of indexing. Is it best to use only one of the index variables available, or to combine several of them, and if more than one variable is used how is the combination to be determined? These questions are of most vital importance in mental measurement, arising in one form or another whenever a new test is devised and standardized. It may be pointed out at once that this study does not attempt any general solutions for these problems, but by analyzing a definite type of test material, aims to carry the investigation a little farther than is possible under incidental treatment in the construction of a particular scale.

Intelligence test material was chosen for two reasons. First, a large body of such data was available. Dr. F.S. Breed and Dr. E.R. Broslich made available their excellent data for three intelligence tests given at two grade levels in The University of Chicago High School. Mr. Guy Capps also turned over two thousand copies of the Terman Group Test of Mental Ability, Forms A and B, administered and repeated in a number of Kansas high schools. A second reason for using intelligence tests in this study was that the statistical constants resulting from such material are more stable than from any similar test data of which the writer is aware. The violent and often inexplic-

able fluctuations in such constants as correlation coefficients and standard deviations for large, well-known achievement tests for middle and high schools, would only carry to further the analysis. It is a study of this kind, that of historical stability, if not heredity, is due for the most part to the hands of the past and to their careful construction.

Section 3. Experimental Methods

3. The Group Studies

For an analytical study of this type it is desirable that the group studies should possess the following properties. They should be large enough to insure reasonable stability in the statistical constants; and second, they should be as homogeneous as possible. The selection of pupils from large schools is an extremely difficult to obtain in best work but is essential to the success of the study. Nevertheless, for the purpose described below a group made up of 100 pupils of 1911 was chosen by combining three schools with an enrollment of less than 100 each, and the total group.

The group described here was made up of 100 pupils of 1911 from the University of Chicago Laboratory School. Groups I, II, and III, 33, 33, and 34 pupils respectively, were respectively from the University of Chicago Laboratory School. The 7th Grade pupils were chosen from the 7th grade list, and were prepared to enter the regular first grade of the following year, there being no eighth grade in the Laboratory Schools.

Group I High B differed from I High only by the addition of 10 pupils in certain of the tabulations. While the groups described are undoubtedly select, they are unusually homogeneous as regards social status, training, experience with tests, and age. The largest group, I High C, consisted of 135 pupils from the Rolla, Salem, and St. James public high schools of Kansas, selected as described above.

The age distributions for these groups are given in Table 3. The Kansas high school pupils were a year older than those in The University High School, while the latter were about a year and a half older than the 7th grade group. The distributions in each case present a fair degree of symmetry, the standard deviations increasing with the size of the group.

TABLE 3.-AGE DISTRIBUTIONS FOR THE GROUPS STUDIED

Age	Grade 7	I High A	I High B	I High C
19.0-19.5	1
18.5-18.9	1
18.0-18.4	1	2
17.5-17.9	6
17.0-17.4	1	6
16.5-16.9	..	2	2	13
16.0-16.4	..	3	2	13
15.5-15.9	..	3	3	17
15.0-15.4	..	6	7	18
14.5-14.9	2	3	10	21
14.0-14.4	2	11	13	16
13.5-13.9	6	10	13	9
13.0-13.4	7	4	8	6
12.5-12.9	16	..	3	3
12.0-12.4	10	2	1	..
11.5-11.9	6	1
11.0-11.4	1
Total	50	50	60	135
Mean	12.8	14.1	14.4	15.4
S.D.	.8	1.0	1.1	1.4

b. The Tests Used

Three intelligence tests were administered to the Laboratory School groups by Mr. Breslich during the year 1920-21. The scales used were the Otis Group Intelligence Scale, Form A, the Terman Group Test of Mental Ability, Form B, and the Chicago Intelligence Scale, Form B. These tests were all carefully scored by the writer for Attempts, Rights, ~~and~~ Wrongs, and Accuracy, as well as for the author's score.

The data from the Kansas schools consisted of the Terman Group Test of Mental Ability Form B, and the same tests Form B given the next day. From this group, therefore, it was possible to obtain the reliability coefficients. The large amount of labor involved in scoring each paper for author's score, total Attempts, Rights, Errors, and Accuracy, and carefully checking all of the work is largely responsible for limiting the sample to 135 cases when over 1000 were available.

In addition to the above data, school marks were obtained for 39 pupils in Grade 7. Yearly grades were obtained for Mathematics, English, and History. These marks were converted into a rough scale dividing pupils into seven categories for purposes of correlation. It is the belief of the writer that this degree of accuracy in treatment is about all such data warrent, inasmuch as these marks are intended to give only a rough estimate of the

achievement in the various subjects.

Section 3 Test Administration and Scoring

a. Primary and Secondary Index Variables

In administering and scoring a test the following variables must of necessity be taken into consideration directly or indirectly: Difficulty, Time, Attempts, Rights, Wrongs, and Omissions. This implies of course that the test material consists of a series of items where the responses may be scored right, wrong, or omitted. The six variables listed above will be referred to hereafter as Primary Index Variables and any function involving more than one of them as secondary index variables. Thus if Accuracy be defined as Rights divided by Attempts, such a score would be termed a Secondary Index Variable. Also in order to save space these variables will usually be denoted by the initial letter of each word, i. e.,

D = Difficulty
A = Attempts
T = Time
R = Rights
W = Wrongs
O = Omissions

b. The Difficulty Factor Eliminated

The problem of the difficulty of the various items in the tests is one which must be settled before proceeding farther. This problem is implicitly solved by the authors of the tests wherein each item is given an equal or point

value with all others. Two questions arise in this connection: are the items of equal difficulty, and if not should they be weighted to obtain an accurate score? It is a well known principle in the theory of index numbers that the longer the series the less the effect of differences in the weights of the individual items. Test scores are, of course, really index numbers. It may be readily conceded therefore, that the authors of intelligence scales consisting of so long a series of items, are quite justified in assigning equal weights to each item regardless of the better values which might be assigned on theoretical grounds. By way of justifying this assumption, an example will be given of a short series with considerable variation in weight from item to item. If differences in weight are not significant in so short a series, they will be even less so for a very long one. This method is that of the unfavorable case.

Test 2 of the Terman Group Intelligence Test consists of 11 items the response to each being a best answer appropriately checked. One hundred pupils were selected who had taken both Forms A and B a day apart. The test papers for Form A were then scored for errors and rough weights assigned in the usual way assuming a normal distribution of difficulty.

TABLE 3.--PER CENT FAILING AND WEIGHTS FOR TERMAN TEST 2

Item	1	2	3	4	5	6	7	8	9	10	11
Per cent Failing	3	19	22	30	36	24	16	47	65	32	55
Weight	2	3	4	4	4	4	3	5	7	4	5

Next by exceedingly tedious computation a weighted and an unweighted score for each pupil was obtained. Thus a boy responding correctly to items 1, 3, 7, 9, and 11 on the test received a weighted score of 31 and an unweighted score of 5. the number right. A correlation table was then made for weighted and unweighted score with a resulting coefficient of $R_{uw} = .972$. The standard error in estimating unweighted from weighted score, $S_e = .9$. The correlation between the two unweighted forms of the test was then obtained, giving $R_{AB} = .597$. The standard error S_e in estimating unweighted form A from unweighted B was 4.8.

Thus the correlation between two forms of the same test is much lower than that between weighted and unweighted scores on the same form. Also the standard error of estimate of unweighted from weighted is about one-fifth that from form B. The weighting of the items then gives a degree of refinement considerably beyond the reliability of the test itself i.e. correlation of two forms. For 40 weighted and unweighted items of different material the writer obtained correlations of .998, .997, .995 with the corresponding reliability coefficients between .85 and .90. All of these results point to the conclusion that for a fairly long series weighting of the separate items is unnecessary. A complete solution of the problem would involve experimentation with series of various lengths, items of various difficulty, and populations of different size. Such exhaustive treatment is clearly beyond the scope of this study. It may be finally pointed out that a number of recent achievement tests have appeared first with weighted items and later with weights dropped when it

was realized that a slight difference these made in the resulting scores.

C. Methods of Administering the Tests

With the factor of difficulty eliminated the scoring or indexing problem is greatly simplified. The plan may now be described as the method of Unit Responses, the response to each item being scored as a unit point. Furthermore if omissions be neglected or counted as errors, the remaining primary index variables are reduced to Time, Attempts, Rights, and Wrongs with the relationship,

$$A = R+W$$

if omissions be counted as errors. This assumption will be made in a subsequent discussion. The number of omissions occurring in the tests used was negligible. It thus appears that four variables, T, A, R, and W will have to be studied, the last three not being independent. Also all scoring formulae or secondary index variables will be functions of these four primary indexes.

Two plans for administering such tests are possible. One plan is to fix the number of Attempts allowing Time, Rights, and Errors to vary, while the second method is to fix the Time, allowing Attempts, Rights, and Wrongs to vary. In the last analysis then, only three index variables need to be considered, the fourth being arbitrarily fixed by the plan of administering the tests. According to the first scheme outlined, one allows all of the children to

THE UNIVERSITY OF CHICAGO
DIVISION OF THE PHYSICAL SCIENCES
DEPARTMENT OF CHEMISTRY

REPORT OF THE
COMMISSIONERS OF THE
UNIVERSITY OF CHICAGO
FOR THE YEAR 1900
PUBLISHED BY THE
UNIVERSITY OF CHICAGO
PRESS

CHICAGO
1901

finish the test thus keeping Attempts constant. The time is then recorded by stop watches or clock device and Rights and Wrongs obtained from the papers. By the second plan a fixed time limit is set for all the pupils, Attempts, Rights, and Wrongs being then scored on the test papers. The second method is obviously simpler than the first, and is now followed in the great majority of tests of all kinds. Certain tests, it is true, neglect actual time but these are not considered here. The material used in this study is all administered under the plan of fixing time giving the three primary index variables A, R, and W.

d. Authors' Plans of Scoring

For the three intelligence tests described above, the authors have set forth scoring formulae for each test of the battery. These formulae are obviously expressed as functions of the primary variables A, R, and W, while relationship $A = R + W$ makes it possible to set down the equations in terms of any two of the variables. In the following tabular scheme, therefore, R and W have been employed throughout.

TABLE 4.- AUTHORS' SCORING FORMULAE ON THE THREE SCALES

Scale	Test									
	1	2	3	4	5	6	7	8	9	10
Otis	R	R	R-W	R	R	R	R	R	R	R
Terman	R	2R	R-W	R	2R	R-W	R	R-W	R	R
Chicago	$R - \frac{1}{2}W$	2R	$2(R - \frac{1}{2}W)$	R	2R	—	—	—	—	—

Figure 1. The effect of the concentration of the Fe^{2+} on the rate of the reaction of Fe^{2+} with H_2O_2 at $[\text{H}_2\text{O}_2] = 0.001 \text{ M}$, $[\text{Fe}^{2+}] = 0.001 \text{ M}$, $[\text{H}^+] = 0.01 \text{ M}$, $[\text{H}_2\text{O}] = 55.5 \text{ M}$, $T = 25^\circ\text{C}$.

1. The first step is to identify the problem or question that needs to be answered. This involves understanding the context and the specific requirements of the task.

1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2632, 2633, 2634, 2635, 2636, 2637, 2638, 2639, 2640, 2641, 2642, 2643, 2644, 2645, 2646, 2647, 2648, 2649, 2650, 2651, 2652, 2653, 2654, 2655, 2656, 2657, 2658, 2659, 2660, 2661, 2662, 2663, 2664, 2665, 2666, 2667, 2668, 2669, 2670, 2671, 2672, 2673, 26

7

It will be seen from table 4 that two of the scales consist of ten tests each while the third includes only five. Considerable variation in the scoring formulae is also to be noted. The Otis Scale has 2 tests scored by R and only one by R-W; Ternan, on the other hand, scores four tests by R, 3 by 2R, and 3 by R-W. For the Chicago Scale two new formulae appear, $R - \frac{1}{2} W$ and $2(R - \frac{1}{2} W)$. All of the above formulae are clearly special instances of the general linear form

$$S = a(R + bW)$$

where a and b are constants.

The simplicity of the Otis' scoring formulae immediately raises the question as to the advisability of scoring the other scales by the same method. This problem is discussed in the following sections. Another question concerns the formulae R-W and $R - \frac{1}{2} W$. These formulae are employed for material of the True-False type and three-choice variety, on the supposition that they correct for the element of guessing for such tests. This problem will also be taken up in some detail in later discussion. Finally the doubling of scores for individual tests is a point that needs consideration. So far as the writer is able to determine, this doubling was effected because the author of the test felt such tests to be worth about twice as much as others, or wanted to increase the total points possible to some convenient number. The point at issue is the same as that

between weighted and unweighted items. A graduate student made a study of the Chicago Scale regarding this problem and found that the weighting of the three tests as indicated in Table 4 affected the correlation but slightly, the coefficient between weighted and unweighted scores being .91.

It thus appears that Rights are the basis of most of the scoring formulas employed on these tests and that other forms have been used to correct Right responses for guessing or to weigh the whole test because of the relative importance in the battery making up the scale.

Section 4. Methods Employed

The general method of this study will first be to analyze the interrelationships of the index variables involved, and then to set up certain criteria of good indexes and attempt to evaluate the variables in terms of these. Obviously this is an indirect approach and it must necessarily be so from the nature of the problem. The actual technique employed will involve a considerable amount of correlation, for this is the best method of studying the relationships between index variables from tests. For the case of the cubes described above the index variables

were functionally related, i.e.

$$\text{volume} = (\text{edge})^3$$

Such functionality can only be approached by empirical data, the correlation coefficient giving the most convenient approximation for linear functions. Thus if score and Right are correlated to the extent of .93 with regression linear, there is a very close approximation to the functional relationship $S = KR$. The chief advantage of correlation is that it gives a numerical estimate of the closeness of such relationships or approach to linear functionality.

Index variables will be analyzed by batteries and by single tests. Age, school marks, and other intelligence tests will be used as criteria against which to check the various formulae.

All of the calculations below have been performed by the writer and have been checked with care. Correlations were obtained by the usual product moment method with a specially designed correlation form. This was found to cut down the labor of calculation very materially especially when "batteries" of coefficients were required. Blakeman's test for linearity was applied in a few cases with the result that the writer believes the great majority of the tables exhibited sufficient linearity so as not to reduce the correlation beyond the limits of probable error.

Part I will be concerned with an analysis of the five index variables by whole scales. By the indirect method of correlation, the relative merit of the indexes will be determined. Part II will involve a similar analysis of the individual tests of the scales. In Part III certain form-

ulas will be developed and their validity and reliability determined in mathematical terms.

Part I Analysis of the Index Variables by Whole Scales

A.- Comparative Validity and Reliability of the Indexes

The most direct approach to the problem of indexing will be to compare the author's plan of scoring as exhibited in Table 1 with the similar plans of counting merely Attempts, Blasts, Errors, and Determining Accuracy for whole scales. This procedure will reveal the relative merit in reliability of these various primary variables when several tests are needed to give a scale score. The term "Scale Score" is here employed to distinguish it from test score which will be used to denote the score on one of the components making up the total scale. Thus the Otis Scale is made up of 16 component tests.

Five types of comparison will be made in evaluating the index variables for general reliability. These include:

- a. Relationship between index variables on the same scale
- b. Relationship between scales by means of index variables
- c. Reliability of a scale by different index variables
- d. Correlations with age
- e. Correlations with school marks

Section B Relationship between Factors Computed on the Same Scales

As pointed out above the reliability of each of the var-

ibles cannot in general be functional for empirical data such as those obtained from mental tests. Nevertheless it will be valuable to discover the closeness of linear relationships as indicated by the correlation coefficient.

These results are set forth in Table 5. All of the correlations were computed on total Attempts, Rights, etc. for the entire series. The agreement from group to group and scale to scale between correlations for the same two variables is striking and may be viewed with pardonable satisfaction by one who has viewed many ineliminable differences for other types of tests and smaller groups.

TABLE 5.- CORRELATIONS BETWEEN VARIABLES FOR THE SAME SCALES

Scale and Group	Pairs of Variables Correlated									
	S×A	S×R	S×W	S× $\frac{R}{A}$	A×R	A×W	A× $\frac{R}{A}$	R×W	R× $\frac{R}{A}$	W× $\frac{R}{A}$
Otis 7	+.73	+.22	-.53	+.73	+.73	+.21	+.13	-.51	+.75	-.21
Otis I B	+.67	+.31	-.59	+.61	+.37	+.30	+.01	-.37	+.77	-.31
Ter. 7	+.52	+.26	-.45	+.66	+.70	+.41	-.22	-.30	+.30	-.34
Ter. I B	+.52	+.30	-.67	+.33	+.71	+.21	+.11	-.51	+.73	-.21
Chi. 7	+.35	+.52	-.54	+.60	+.62	+.21	-.10	-.32	+.32	-.22
Chi. I B	+.52	+.37	-.71	+.35	+.63	+.12	+.00	-.65	+.77	-.94
Mean	+.62	+.35	-.58	+.70	+.61	+.31	-.00	-.21	+.62	-.38

for example the correlation between authors' score and total Attempts for Otis Grade 7 is $.73 \pm .05$, while for Otis I High B it is $.67 \pm .05$. By the formula $PE_{a-b} = \sqrt{(PE_a)^2 + (PE_b)^2}$ the difference between the two correlations may be written in the form $\text{diff.} = .05 \pm .07$, indicating that a difference of .05 is insignificant. To be significant such a difference would have to be 3 times its P.E. Similar comparisons between groups show that nearly all diff-

ences may be accounted for by the fluctuations in sampling. Moreover if tests made on the same group with different scales be considered as independent samples, most of the inter-scale differences are also insignificant.

The table yields some very interesting results. In the second column it will be observed that the correlation between S and R are very high, the mean of the six coefficients being $+.96$. The extremely high correlation of $.99$ for the Otis scale is due to the fact that only one of the tests in this scale is not scored directly by Rights. The above result, therefore, raises a question as to the advisability of scoring this lone test differently from the rest. Moreover the two correlations of $.96$ for the Terman Scale indicate that even with a fairly complicated system of scoring, the agreement between authors' score and total Rights is extremely close. With still more complicated formulae in the Chicago Scale the correlations between S and R are again very high. For whole scales, then, a mere enumeration of the total number of right responses gives a result very nearly proportional to that obtained by the use of various formulae for the individual tests making up the scale. It is to be noted, however, that the above result appears to be valid when single tests are pooled to give a total score. For single tests, changes in the scoring formulae have marked effect upon the correlations with criteria as will be shown later.

Accuracy furnishes the next highest correlation with Score, the mean coefficient being +.76, while A and W come next with correlations of +.59 and -.50 respectively. If Score be adopted as a criterion, therefore, the best index variables in order would be Rights, Accuracy, and Attempts or Errors. It is rather surprising that merely counting the Attempts or Errors on these tests gives so good an index of intelligence.

In the last column of the table a very high and consistent negative correlation is found between W and $\frac{R}{A}$. This means that a pupil who makes a great many errors is very inaccurate as measured by the index $\frac{R}{A}$, or that such an index variable is an exceedingly good measure of what is ordinarily understood by Accuracy. It is of some interest to note that the above correlation becomes -1.00 when the number of attempts is fixed and the relationship

$$A = R + W$$

still obtains. In this case the functional relationship

$$R = -W + \text{const}$$

holds strictly, and this implies perfect negative correlation as is shown more fully in a section below.

If the index variable A be used as a measure of speed a number of interesting relationships are brought out by the remaining coefficients. Speed and accuracy are evidently uncorrelated, the highest correlation between these two variables being $-.16 \pm .62$ which is insignificant because it is not even twice its probable error in amount. Moreover

the differences in sign are such as to give a near correlation for the column of less than .51. It will be further noted that R is correlated with A to the extent of +.59 on the average, while W with A gives a near correlation of +.84. An average correlation of +.59 between S and A was previously noted. These results indicate that the way to get a high intelligence score is to work fast. By working fast one is likely to make more mistakes, but he is much more likely to get more items right and make a higher score than if he worked more slowly. Intelligence tests have frequently been called "alertness" tests. The above findings indicate that with considerable error-proneness they might also be termed "speed" tests. The assumption thus far in determining general reliability has been that the authors' score is the best index of intelligence. If, on the other hand, Accuracy had been assumed to be the best index, the speed factor would have been eliminated, there being no general tendency as indicated by the zero correlation for a pupil to get a high or low accuracy score by changing his speed. It will also be recalled that Accuracy was highly correlated with Score (+.76) so that $\frac{R}{A}$ as an index variable is ~~fairly~~ ^{fairly} consistent with S without being unduly influenced as is the latter variable by the undesirable speed factor.

The average correlation of -.49 between R and W means that the more items the pupil gets right, the fewer he is

likely to get wrong. If Attempts are constant the above correlation becomes -1.00 as in the case of Wrongs and Accuracy. Finally the variable R and $\frac{R}{A}$ are found to correlate on the average $+1.00$, a coefficient which becomes $+1.00$ when Attempts are constant. (Theorem 3 Appendix)

Section 6. Relationships between Scales by the Same Index Variable

The relationships between index variables on the same scale have been discussed in the preceding section, with the result that the general order of merit with score as a criterion is Rights, Accuracy, and Attempts or Errors. The variable R possesses a simplicity which considered in connection with its close agreement with S suggests that it might well be substituted for that variable in indexing by batteries of tests. The index $\frac{R}{A}$ was found to be in close agreement with S and R and to possess the advantage of being unaffected by speed.

The next procedure will be to evaluate the index variables indirectly by determining the correlations between pairs of scales indexed by the same variables. The closeness of this correspondence will give a measure of the effectiveness of the particular mode of indexing. The former comparison was inter-variable, the present is inter-test. Table C gives these results for groups I High B and Grade 7. Comparison of this table with Table 5 reveals the fact that the correlations in the former are much more nearly the same size. All of the coefficients are significant, while there are few of the differences between any

two which may not be attributed to sampling by the usual method. The means for the columns and rows of the table bring out the stability more clearly inasmuch as the deviations from these are in most cases very slight.

The really surprising results exhibited by this table are indicated in the column of means for the various indexes. Here it appears that with inter-scale correlation as a criterion, S, R, $\frac{R}{A}$, and W are all about equally good for purposes of indexing, while A is somewhat poorer than the rest. To discover almost as high a correlation between two intelligence scales by merely tabulating total errors as by using the author's score or total Rights is at first somewhat surprising. It is a closer agreement than might be expected from the average correlation of $-.58$ between S and W from Table 5. Inter-test correlation gives a measure of the extent to which two tests agree in measuring the same characteristic. From Table 6 it appears that this agreement is about equally close when any of the four index variables is employed, the general order of merit being indicated in the table as S, R, $\frac{R}{A}$, W, and A. This order, it will be recalled, is in harmony with that found in the preceding section. It may finally be pointed out that the lack of variability among the coefficients strongly suggests that combining these index variables as in the case of S will give but slightly better index than batteries of tests such as these are employed.

TABLE 6.--INTER-SCALE CORRELATIONS BY FIVE INDEX VARIABLES

Index Variable	Pairs of Scales and Groups						Mean	Order
	Term. x Otis		Term. x Chic.		Chic. x Otis			
	El.	H.S.	El.	H.S.	El.	H.S.		
Score	+.72	+.85	+.65	+.76	+.75	+.49	+.76	1
Attempts	+.66	+.72	+.57	+.46	+.57	+.32	+.56	2
Rights	+.73	+.82	+.59	+.71	+.73	+.38	+.74	3
Wrong	+.72	+.83	+.66	+.66	+.66	+.35	+.69	4
Accuracy	+.70	+.98	+.61	+.77	+.66	+.74	+.73	5
Mean	+.71	+.82	+.60	+.63	+.67	+.42	+.70	

Section 7 The Reliability of a Scale by Different Index Variables

Another method of studying the reliability of the different index variables is to obtain the reliability coefficients for a scale under each of the indexes. Group I High C was used for this purpose. It will be recalled that this group consisted of 155 pupils who took Form B of the German Scale and Form A of the same test on the following day. The reliability coefficients in this case are given by the correlations between the variables on the two forms of the scale. These correlations are given in Table 7 with the contingency tables presented in the appendix.

[illegible]

TABLE 7.- RELIABILITY COEFFICIENTS FOR FIVE TESTS
A AND B ON GROUP I HIGH C

Variables	Reliability Coefficient
Score	$+.858 \pm .010$
Attempts	$+.804 \pm .009$
Flights	$+.826 \pm .011$
Wings	$+.787 \pm .011$
Accuracy	$+.841 \pm .017$

It is at once apparent that the subject's score has the highest reliability with a coefficient of .858. A glance at the corresponding correlation table in the appendix will give this result more clearly. Here the linearity of regression is at once apparent, and the remarkable agreement exhibited in graphical form. Close correspondence of this type appears to the writer as one of the most significant achievements of standardized tests. Just what the tests do measure or index is often ambiguous, but to index any mental characteristic with such a high degree of reliability is in itself a most noteworthy achievement. It should also be borne in mind that all such reliability coefficients (and indeed all such correlations) depend upon the group. Selection will in general tend to reduce such correlations while heterogeneity due to such factors as age will tend to increase it.

Returning to the remaining coefficients in the above table, we find that the order of reliability for the five variables is, S, E, A, W, and A. This is precisely the order obtained by the method of inter-test correlation as shown in Table 6.

If the symbol r_{xx} be employed to denote the reliability coefficient for a test indexed by the variable X, the differences between the correlations in Table 7 may be exhibited as follows:

$r_{ss} - r_{AA} =$	$+.334 \pm .034$	difference is significant
$r_{ss} - r_{RR} =$	$+.012 \pm .015$	difference is insignificant
$r_{ss} - r_{ww} =$	$+.171 \pm .029$	difference is significant
$r_{ss} - r_{\frac{RA}{AA}} =$	$+.067 \pm .020$	difference is significant

The formula used for calculating the probable errors of the differences is,

$$P.E_{a-b} = \sqrt{(P.E_a)^2 + (P.E_b)^2}$$

The small differences between the reliability coefficients for S and R is insignificant, while the remaining differences are sufficiently large in comparison with their probable errors to be significant. Thus from the standpoint of reliability the Terman Scale is indexed equally well by author's score or total rights, and next best by accuracy. A reliability coefficient of .74 is generally considered high, and it is remarkable that total errors on the two forms should correlate to such an extent. The difference between r_{ss} and r_{ww} , however, is over five times its probable error so that the reliability of errors as an index is significantly less than for score. Similarly Attempts furnish a much less reliable index than Score, Rights, or Accuracy.

In this part of the study no attempt is made to analyze the individual tests making up the scales. Nevertheless it is interesting to note from Table 4 that three

of the Torrance tests are scored R-W, three 1. and four by R alone. The difference in the reliability coefficient R_{SS} under this plan, and R_{RR} by counting merely rights is +.016, .015 as shown above. Formulas of the type 3R and R-W, thus appear to have no effect on the reliability of the total score, and their use for sub scales is open to question.

Section 3. Correlations with Age

The age factor is always of interest when studying test results. Table 4 which gives the age distributions in half-year intervals, shows ranges for the various groups from four to eight years. In Table 5, the correlations between age and the four index variables are given for Grade 1 and I (Fig. 3). A consistency in these coefficients is at once apparent. The mean correlation of $-.41$ between age and score indicates that the younger pupils are brighter than the older ones within the same grade group. Similarly the mean correlations $-.20$, $-.29$, $-.36$, and $+.37$ show that the younger pupils get more items right, are more accurate, are speedier, and make fewer errors than the older children.

TABLE 8.-- CORRELATIONS BETWEEN AGE AND INDEX VARIABLES

Scale and Group	Age with				
	Score	Attempts	Rights	Wrongs	Accuracy
Otis 7	-.41	-.37	-.37	+.06	-.24
Otis IB	-.40	-.37	-.39	+.27	-.32
Terman 7	-.47	-.22	-.40	+.17	-.37
Terman IB	-.39	-.21	-.37	+.28	-.34
Chicago 7	-.39	-.21	-.35	+.20	-.33
Chicago IB	-.39	-.30	-.39	+.34	-.26
Mean	-.41	-.26	-.38	+.20	-.29

The superiority of the younger child is evident, therefore, no matter which index is employed for the tests.

When arranged according to the size of the correlation with age the order of the index variables appears from the means as, S, R, ^AR, A, and W. The three variables Score, Rights, and Accuracy retain the order found in the previous sections. Moreover the mean correlations for score and rights with age are very nearly the same, so that the superiority of S over R as an index is again slight if any. The usual sampling formula reveals no difference in the correlations that is of statistical significance.

Similar correlations are given in Table 9 for the largest group. The coefficients, though somewhat smaller are in harmony with those of the preceding table. The decrease in size is probably due to the greater range in age.

TABLE 9.- CORRELATIONS OF INDEX VARIABLES WITH AGE FOR
TERMAN SCALE FORM B

Group	Age with				
	Score	Attempts	Rights	Wrongs	Accuracy
Grade 7	-.47+.07	-.32+.09	-.46+.08	+.17+.09	-.37+.09
I High B	-.39+.08	-.31+.09	-.37+.08	+.28+.09	-.34+.08
I High C	-.16+.06	-.14+.06	-.16+.06	+.63+.06	-.18+.06

By increasing the range through several grades the correlation becomes positive. In fact heterogeneity or lack of selection appears to have a curious effect on correlations with age. For children of exactly the same age the coefficient is of course zero; when the range is increased to several years as in the typical grade group, the correlation is negative. As the range in age increases, the coefficient approaches zero again, and finally passes through this value to positive values of considerable size if several grades are pooled to give a long range. The theoretical curve for the correlation coefficient will thus have an appearance resembling that in Figure 4. The negative correlation for the age interval OA has been accounted for by the appearance of older retarded children in grade groups. This explanation, while plausible, does not seem satisfactory for groups such as Grade 7, which is usually free from children of this type. The positive correlation increases from A as the age span is lengthened.

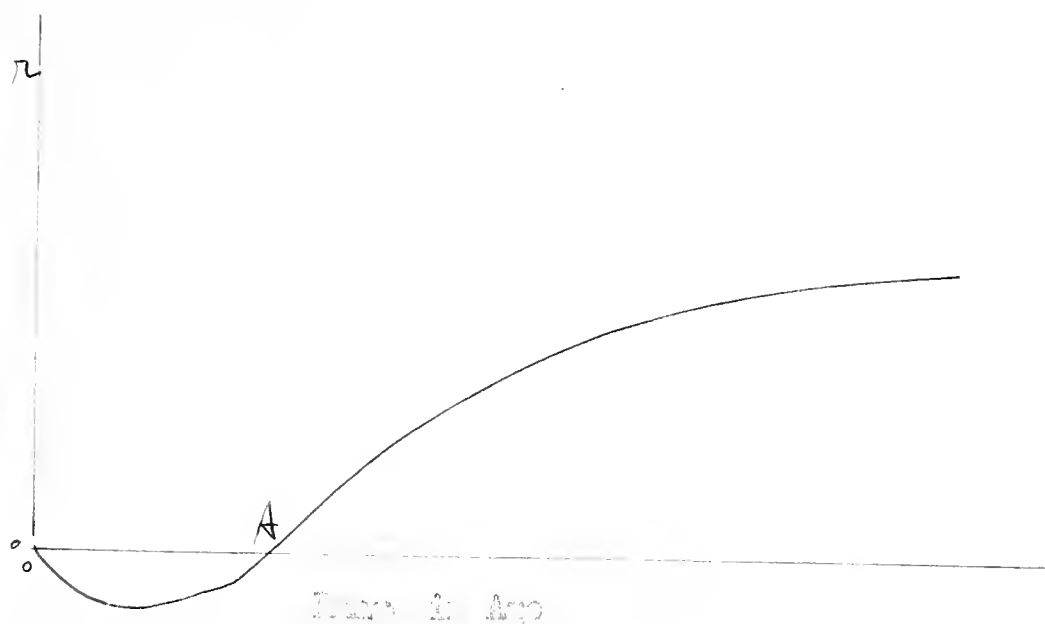


Fig. 4. Theoretical curve for correlation with age

Section 2. Correlations with School Marks

While school marks are obviously in exact estimate of the ability or achievement of pupils, nevertheless they may serve as a useful device and as a means to an end test result. Their relative inaccuracy has been greatly exaggerated. Bart, Brooker, Kelly, and others have shown that the predictive value of such marks is often as high as

that from intelligence or achievement tests. Moreover in the present comparisons the question of accuracy is not of great importance inasmuch as each index is checked against the same set of scales. Whatever unreliability exists in the marks, therefore, will affect all correlations alike.

TABLE 10.- CORRELATIONS BETWEEN MARKS IN ENGLISH AND THE INDEX VARIABLES ON THE THREE SCALES

Scale	Marks in English with				
	Score	Attempts	Rights	Wrong	Accuracy
Otis	+.56	+.36	+.57	-.49	+.46
Terman	+.60	+.37	+.56	-.51	+.46
Chicago	+.58	+.37	+.53	-.53	+.43
Mean	+.57	+.31	+.53	-.51	+.47

TABLE 11.- CORRELATIONS BETWEEN MARKS IN HISTORY AND THE INDEX VARIABLES ON THE THREE SCALES

Scale	Marks in History with				
	Score	Attempts	Rights	Wrong	Accuracy
Otis	+.30	+.17	+.54	-.45	+.55
Terman	+.63	+.34	+.60	-.38	+.53
Chicago	+.48	+.17	+.41	-.57	+.47
Mean	+.55	+.17	+.53	-.46	+.52

TABLE 12.- CORRELATIONS BETWEEN MARKS IN MATHEMATICS AND INDEX VARIABLES ON THE THREE SCALES

Scale	Marks in Mathematics with				
	Score	Attempts	Rights	Wrong	Accuracy
Otis	+.60	+.45	+.61	-.34	+.40
Terman	+.54	+.33	+.77	-.56	+.36
Chicago	+.47	+.39	+.35	-.27	+.37
Mean	+.54	+.35	+.53	-.35	+.34

Correlations of the five index variables with English, History, and Mathematics are given in Tables 10, 11, and 13 respectively. Inspection of the average correlations for these tables shows a close agreement for the three subjects studied. The coefficients for English are highest, History next, and Mathematics lowest, but the differences are slight. In all three tables Score has the highest correlation with school marks. The next highest correlations in order are Rights, Accuracy, Wrong, and Attempts. This is precisely the order found in the section on reliability. Thus if school work be measured by marks the relative merit of the various indexes for prediction is $S, R, \frac{R}{A}, W$, and A . As in the preceding sections, correlation involving S and R are more nearly equal than the others.

Section 10. Summary for Reliability of Indexes for Whole Scales

For whole scales consisting of batteries of tests, the authors' formulae appear to be slightly superior to total Rights as an index. Table 13 gives the average correlations and differences in favor of S (absolute values considered)

TABLE 13.- SUMMARY OF CORRELATIONS FOR SCORE AND RIGHTS

Variables Correlated	Average Coefficient	Difference in Favor of S
Score and Rights	+0.16	
Scales Indexed by Score	+0.14	
Scales Indexed by Rights	+0.1405
Terman Forms B and A by S	+0.14	
Terman Forms B and A by R	+0.1001
Age with Score	+0.14	
Age with Rights	+0.0803
Marks with Score	+0.55	
Marks with Rights	+0.5503
Total Difference		.08

The extreme simplicity of scoring by rights, however, would seem to more than outweigh the slight advantage in favor of more complicated formulas.

Accuracy has been shown to have the peculiar advantage of being unaffected by speed, and at the same time to possess high reliability. The summary correlations in favor of Score are shown in Table 14. The total differ-

TABLE 14.- SUMMARY OF CORRELATIONS FOR SCORE AND ACCURACY

Variables Correlated	Average Coefficient	Difference in Favor of S
Score and Accuracy	+0.76	
Scales Indexed by Score	+0.76	
Scales Indexed by Accuracy	+0.7603
Terman Forms B and A by S	+0.71	
Terman Forms B and A by R	+0.6607
Age with Score	+0.71	
Age with Accuracy	+0.6612
Marks with Score	+0.75	
Marks with Accuracy	+0.7411
Total Difference		.53

ences in favor of S are shown in Table 14. The total differences in favor of S indicates that Accuracy is somewhat less satisfactory than R according to the criteria employed. Moreover it is a more involved complicated index than R, but not so involved as S.

The results for Errors are presented in Table 15. The general merit of Wrong as an index is less than that of the preceding variables. Errors, however, have a surprisingly high reliability and are utilized to advantage in formulas discussed in the following sections.

TABLE 15.- SUMMARY OF CORRELATIONS FOR SCORE AND ERRORS

Variables Correlated	Average Coefficient	Differences in Favor of Score
Score and Wrong	-.50	
Scales Indexed by Score	+.76	
Scales Indexed by Wrongs	+.6067
Terman Forms B and A by S	+.31	
Terman Forms B and A by W	+.7217
Age With Score	-.41	
Age With Wrongs	+.3031
Marks With Score	+.55	
Marks With Wrongs	-.3335
Total Difference		.68

Attempts, which are frequently used as an index for tests, appear to have the least merit of any of the variables discussed. Table 16 gives the averages and differences as in the above tables. The total absolute difference in favor of S is greater than for any of the preceding variables.

TABLE 16.-- SUMMARY OF CORRELATIONS FOR SCORE AND ATTEMPTS

Variables Correlated	Average Coefficient	Difference in Favor of Score
Score and Attempts	+.50	
Scales Indexed by Score	+.76	
Scales Indexed by Attempts	+.5620
Terman Forms B and A by S	+.91	
Terman Forms B and A by A	+.6823
Age with Score	-.41	
Age with Attempts	-.2615
Marks with Score	+.55	
Marks with Attempts	+.3827
Total Difference		.85

B.-- The Discriminative Capacity of the Indexes

In addition to the general reliability of an index, another valuable property of such a variable is the extent to which it makes possible discrimination between individuals and between groups when real differences exist. A test which reveals too narrow a range for a given group fails to discriminate between the individuals of that group. Such undistributed score is a defect in the test or in the mode of indexing. Similarly a test or mode of indexing which fails to discriminate between groups is defective if the characteristic is in reality different in type from group to group. Thus a test which shows all individuals in Grade 5 to possess the same ability, and at the same time reveals

no difference between mean scores for Grade 5 and 6 is lacking in individual and in group discrimination. The fundamental assumption is, of course, that such individuals and groups do vary and that failure to detect the variations lies in the particular mode of indexing the trait in question.

Section 11 Capacity of the Indexes to Discriminate between Individuals

Discrimination between individuals of a group is best studied by means of frequency distributions. In the present study, however, such an elaborate method as this is unnecessary inasmuch as intelligence tests are of sufficient length to give a fairly good spread for all indexes. The distributions for S, R, A, W, and \bar{R}_A in the Appendix are typical of those for all three intelligence scales. The standard deviation for these variables are given in Table 17.

TABLE 17.- STANDARD DEVIATIONS FOR TERMAN FORMS A AND B WITH GROUP I HIGH C

Variable	\bar{C}_B (first)	\bar{C}_A (second)	$\bar{C}_B - \bar{C}_A$	Diff. $\bar{P}.E. diff.$
Score	35.74 \pm 1.39	30.31 \pm 1.97	5.43 \pm 1.36	1.5
Attempts	24.14 \pm 0.99	20.63 \pm 0.65	3.51 \pm 1.30	2.7
Rights	36.95 \pm 1.11	35.54 \pm 0.97	1.41 \pm 1.47	2.3
Wrongs	19.73 \pm 0.31	19.63 \pm 0.81	0.10 \pm 1.15	0.1
Accuracy	0.12 \pm 0.01	0.18 \pm 0.01	0.06 \pm 0.01	1.0

There is some evidence that there is less variability in performance on the second trial (Form A) than on the first (Form B). This is incidentally a bit of evidence to the

effect that equal practice for a group of pupils tends to bring them more closely together about a central type, a result contrary to that held by some psychologists. The differences, however, are slight, although in one direction, and the two forms of the test may not be equivalent for this purpose. The result is then merely suggestive.

The standard deviations for A, R, and W in Table 17 admit of direct comparison inasmuch as they are all expressed in point or response units. The order of discriminative capacity for these variables is then R, A, and W. The indexes S and R are expressed in different units and hence may not be compared with the rest. Considered on the point basis however, author's score has the greatest capacity for discrimination between individuals on account of the weighting and formulae involved. The standard deviations for Grade 7 and I High B are also given in Table 18. The results agree with those of the preceding table.

TABLE 18.- STANDARD DEVIATIONS FOR THE THREE SCALES

Scale and Group	Standard Deviations for				
	Score	Attempts	Rights	Wrong	Accuracy
Otis Grade 7	19.91	18.00	19.22	13.48	0.08
Otis I High B	22.40	18.02	21.55	17.02	0.08
Terman Grade 7	23.52	19.17	18.32	12.46	0.08
Terman I High B	26.87	18.39	22.34	14.87	0.09
Chicago Grade 7	10.52	6.20	6.15	4.60	0.09
Chicago I High B	12.03	5.80	8.24	6.50	0.10

In order to study the variability of a group by a statistical measure independent of the units employed, Pearson's Coefficient of Variation, $V = \frac{100 \text{ S.D.}}{M}$, was employed. The results for two groups appear in Table 19. It is at once apparent that while V is independent of the units employed it may nevertheless lead to results which are confusing. The largest coefficients of variation are for W , an index which might readily be supposed to furnish the least variability. The result is brought about by the relatively large standard deviation of W (Table 18) and the relatively low mean (Table 20, below).

TABLE 19.- COEFFICIENTS OF VARIATION FOR THE THREE SCALES

Scale and Group	Coefficients of Variation for				
	Score	Attempts	Rights	Wrong	Accuracy
Otis Grade 7	14.2	10.2	14.3	36.2	9.86
Otis I High B	14.8	9.5	14.2	45.6	10.35
Terman Grade 7	18.8	13.8	16.6	45.5	9.85
Terman I High B	18.7	11.5	17.4	47.7	11.37
Chicago Grade 7	19.5	11.7	14.6	43.0	10.71
Chicago I High B	21.7	10.1	18.8	48.1	13.18
Mean	17.8	11.1	16.0	44.3	10.89

The coefficient of variation, depending as it does upon the position of the distribution on the scale, is likely to give a very misleading result for distributions such as these above, and should in general be avoided for comparisons of this type.

Section 12 Capacity of the Indexes to Discriminate
between Groups

Table 20 gives the means on the three scales for Grade 7 and for I High B. It is at once evident that the second group has the higher mean for nearly all of the indexes. Accuracy, however, appears to be nearly consistent for all three scales and for both groups. From the standpoint of discrimination, therefore, this index is of little value. The correlation tables in the Appendix show a considerable spread for Accuracy while the constants from Table 17, 18, and 19 indicate the extent of this variability.

TABLE 20.- MEANS FOR THE THREE INTELLIGENCE SCALES

Scale and Group	Means for				
	Score	Attempts	Rights	Wrong	Accuracy
Otis Grade 7	139.8	177.0	139.6	37.2	0.80
Otis I High B	151.3	189.3	152.0	37.3	0.80
Terman Grade 7	125.0	138.8	111.4	27.4	0.80
Terman I High B	143.7	159.5	128.3	31.2	0.81
Chicago Grade 7	54.0	52.8	42.1	10.7	0.80
Chicago I High B	55.4	57.2	43.7	13.5	0.76

Individuals within a group, then, differ considerably in accuracy. When the above inter-group comparison is made however, Accuracy is found to be relatively constant. These results indicate that the growth curve for accuracy ordered relative to age will be relatively flat in comparison with ordinary score. This whole matter will be

A	B	C	D	E
F	G	H	I	J
K	L	M	N	O
P	Q	R	S	T
U	V	W	X	Y
Z				

fully treated by the writer in a forthcoming article on Growth Curves under Different Modes of Indexing.

In order to bring out such inter-group differences more clearly they are presented in full in Table 21. The quantity $\frac{D}{P.E.}$ denotes the inter-mean difference divided by the probable error of this difference calculated in the manner explained in preceding sections. Such a quotient gives a convenient index of discrimination. Indexes less than 2 or 3 show that the discriminative capacity of the test for such variables is not significant.

TABLE 21.- DISCRIMINATIVE CAPACITY OF THE INDEXES AS SHOWN BY INTER-MEAN DIFFERENCES IN GRADE 7 AND I HIGH B

Scale	Inter-Mean Difference and Probable Errors for									
	Score		Attempts		Rights		Wrong		Accuracy	
	Diff.	P.E.	Diff.	P.E.	Diff.	P.E.	Diff.	P.E.	Diff.	P.E.
Otis	11.5	2.7	12.3	2.3	12.4	2.7	0.1	2.0	0.0	0.01
Terman	18.7	3.2	20.7	2.4	16.9	2.6	3.8	1.8	0.01	0.01
Chicago	1.4	1.5	4.4	0.8	1.6	0.2	2.8	0.7	-0.04	0.014
Ave. $\frac{D}{P.E.}$	3.7		6.5		4.3		2.1		1.3	

The five variables in order of their capacity are A, R, S, W, and R. Thus the groups studied show the greatest difference with respect to speed and the least with respect to accuracy. This result is quite in agreement with common teaching experience. Pupils can be easily made to hurry, but it is exceedingly difficult to train them to be accurate.

While A shows the best capacity for inter-group discrimination, it is not superior to the other variables for differentiating individuals. Score and Rights again appear to be superior to the other indexes for individual discrimination, a property which is more important than inter-group differentiation.

Section 13 Practice Effect with Repetition

Form B of the Terman Scale was given to group I High C and Form A of the same test given the following day. Assuming that these two forms are equally difficult a practice effect for each of the variables may be noted as in Table 22. There is a positive difference between the means for each of the variables except W. This last negative difference also means an improvement on second trial, so that the practice effect is indicated on all of the variables. The last column shows the significance of this gain. The indexes R, A, S, and $\frac{R}{A}$, reveal gains that almost certainly cannot be accounted for by chance fluctuations, while the change in W is in harmony with that of the other variables. Errors and Accuracy show changes of less significance for practice effect.

TABLE 22.- MEANS FOR TERMAN FORMS A AND B IN GROUP I HIGH C

Variables	M_A (second)	M_B (first)	$M_A - M_B$	Diff. P. E. diff
Score	101.74 \pm 1.72	84.11 \pm 1.36	17.63 \pm 2.35	6.7
Attempts	161.00 \pm 1.20	148.19 \pm 1.40	12.81 \pm 1.84	7.0
Rights	122.41 \pm 1.37	87.15 \pm 1.36	35.26 \pm 2.08	7.3
Wrong	58.59 \pm 1.14	61.04 \pm 1.14	-2.45 \pm 1.60	-1.5
Accuracy	0.68 \pm 0.01	0.59 \pm 0.01	0.09 \pm 0.01	4.4

Part II

Analysis of the Index Variables by Component Tests

The authors' plans of scoring given in Table 3 show that 9 of the tests making up the Otis Scale are scored by the formula $S = R$. These nine tests were therefore chosen for analytical study. The stability of correlations for whole scales has been shown in Part I. In the following sections the intercorrelations of the component tests show a high degree of consistency. The coefficients in general are lower than for whole scales but they indicate the same relationships between index variables. It will also be shown that pooling tests increases both the validity and the reliability of the indexes, an effect which may be roughly forecast by certain predictive formulae.

Section 14 Intercorrelations of Variables for the Otis Components

The correlations between Index Variables for the same components are given in Table 25. In the last line of the table the coefficients for all nine tests pooled are given for comparison.

TABLE 23.- CORRELATIONS BETWEEN INDEX VARIABLES ON NINE OF THE COMPONENTS OF THE OTIS SCALE

Test	Grade 7			I High A		
	A-R	A-W	R-W	A-R	A-W	R-W
1	+.37	+.34	-.81	+.19	+.10	-.77
2	+.33	+.43	-.10	+.31	+.35	-.05
4	+.72	+.50	-.49	+.63	+.66	-.75
5	+.72	+.64	-.66	+.70	+.35	-.63
6	+.58	+.43	-.44	+.47	+.30	-.55
7	+.55	+.56	-.33	+.50	+.35	-.52
8	+.50	+.50	-.33	+.50	+.26	-.86
9	+.31	+.30	-.37	+.81	+.40	-.43
10	+.37	+.50	-.51	+.76	+.32	-.85
Mean	+.61	+.50	-.21	+.50	+.33	-.60
All	+.75	+.51	-.33	+.81	+.31	-.70

It is evident that these are higher than the means of the nine correlations on component tests except for Attempts with Errors, in which case the pool gives the lower value. In certain cases, therefore, pooling or lengthening the tests has the effect of increasing the correlation between the indexes. The exception in this instance is worthy of note as a warning against applying general rules for the correlation on lengthened tests. The high degree of consistency in the coefficients indicates that pooling of such components is a justifiable procedure inasmuch as the test material is fairly homogeneous for purposes of indexing.

Intercorrelations between Rights on the nine component parts of the otis scale are given in Tables 24 and 25 for Grade 7 and I High A respectively. Both groups consisted of 50 pupils. All coefficients larger than three times

their probable errors are printed in heavy type.

TABLE 24.- CORRELATIONS BETWEEN RIGHTS ON THE NINE OTIS COMPONENTS FOR GRADE 7

Test	1	2	4	5	6	7	8	9	10
1		+.54	+.40	+.47	+.50	+.54	+.33	+.50	+.34
2	+.34		+.37	+.52	+.39	+.43	+.26	+.27	+.27
4	+.40	+.37		+.35	+.43	+.37	+.40	+.38	+.19
5	+.47	+.52	+.35		+.51	+.32	+.35	+.14	+.18
6	+.50	+.39	+.43	+.51		+.56	+.32	+.36	+.15
7	+.54	+.43	+.37	+.30	+.50		+.13	+.52	+.35
8	+.33	+.26	+.40	+.35	+.39	+.13		+.19	+.29
9	+.50	+.27	+.38	+.14	+.36	+.52	+.19		+.38
10	+.34	+.27	+.19	+.16	+.15	+.35	+.29	+.33	
Mean	+.41	+.30	+.34	+.37	+.40	+.37	+.39	+.34	+.26

A simple calculation will show that this includes all coefficients numerically greater than .37. In Table 24 only 8 of the 36 correlations are not significant, while in Table 25 the same number occur. Most of these low coefficients are found in the correlations with test 10, the mean value for which is lower than for any other test. This one component then appears to be out of harmony with the rest; i.e. to fail to measure the same thing as the other tests of the battery. Inspection of the Otis Scale shows that test 10 is for memory, a trait quite different from those involved in the other components. Except for this one test a fair degree of consistency is found for the coefficients in both tables. The means for all 36 coefficients in each table are +.35 and +.36 respectively.

TABLE 25.-- CORRELATIONS BETWEEN RIGHTS ON THE NINE OTIS COMPONENTS FOR I HIGH A

Test	1	2	4	5	6	7	8	9	10
1		+.36	+.23	+.45	+.53	+.33	+.50	+.42	+.34
2	+.30		+.41	+.30	+.40	+.35	+.50	+.47	+.21
4	+.35	+.41		+.33	+.24	+.36	+.34	+.38	+.32
5	+.45	+.30	+.25		+.37	+.38	+.44	+.17	+.33
6	+.39	+.42	+.34	+.32		+.48	+.52	+.45	+.17
7	+.33	+.35	+.30	+.30	+.40		+.44	+.36	+.02
8	+.50	+.50	+.34	+.44	+.37	+.44		+.39	+.19
9	+.42	+.47	+.38	+.17	+.45	+.30	+.20		+.13
10	+.34	+.31	+.32	+.33	+.17	+.02	+.10	+.13	
Mean	+.41	+.40	+.33	+.37	+.40	+.35	+.42	+.35	+.21

Tables 26 and 27 show the correlations between errors for the component tests. Only 14 of the 36 coefficients in Table 26 are significant, the mean for the whole table being +.14. Test 10 shows next to the lowest average correlation with the other tests, so that it is of little significance indexed by the rights or wrongs. In Table 27 the coefficients are somewhat higher, the mean of the whole table being +.39. Eight of the 36 correlations are significant, with five of the lowest values appearing with Test 10. It is difficult to explain the difference in correlation for the two groups when indexed by W. Tables 24 and 25 showed mean values nearly identical, but the difference between the mean coefficients for W is too large to be ascribed to chance. One explanation of this difference may be found in the fact that Group I High A made more errors than Grade 7 (See Table). The effect of this was to give less jarring in the contingency tables with a resultant higher correlation.

TABLE 26.- CORRELATIONS BETWEEN ERRORS ON THE NINE OTIS COMPONENTS FOR GRADE 7

Test	1	2	4	5	6	7	8	9	10
1		+.16	+.17	+.49	+.11	+.39	+.30	-.01	+.15
2	+.16		+.24	+.15	+.18	+.35	+.16	+.17	+.13
4	+.17	+.24		+.33	+.30	+.44	+.37	+.33	+.28
5	+.49	+.15	+.32		+.18	+.45	+.35	+.09	+.19
6	+.11	+.18	+.30	+.18		+.40	+.34	+.18	+.22
7	+.39	+.35	+.44	+.45	+.40		+.42	+.31	+.23
8	+.30	+.16	+.37	+.35	+.34	+.43		+.21	+.14
9	-.01	+.17	+.33	+.09	+.18	+.31	+.31		+.12
10	+.15	+.13	+.28	+.19	+.22	+.23	+.14	+.12	
Mean	+.21	+.19	+.31	+.30	+.33	+.35	+.26	+.16	+.18

TABLE 27.- CORRELATIONS BETWEEN ERRORS ON THE NINE OTIS COMPONENTS FOR I HIGH A

Test	1	2	4	5	6	7	8	9	10
1		+.44	+.42	+.38	+.50	+.50	+.47	+.46	+.39
2	+.44		+.44	+.36	+.54	+.57	+.50	+.51	+.18
4	+.43	+.44		+.37	+.37	+.51	+.51	+.39	+.32
5	+.38	+.36	+.30		+.18	+.52	+.40	+.10	+.39
6	+.50	+.54	+.37	+.43		+.48	+.49	+.26	+.16
7	+.50	+.57	+.52	+.52	+.40		+.54	+.50	+.31
8	+.47	+.50	+.51	+.48	+.49	+.54		+.31	+.34
9	+.46	+.51	+.32	+.10	+.30	+.50	+.31		+.10
10	+.39	+.18	+.22	+.29	+.16	+.31	+.14	+.10	
Mean	+.43	+.45	+.40	+.35	+.40	+.51	+.41	+.32	+.24

Comparison of the four tables above shows that the intercorrelation of errors on the 9 component tests is about as high as for Rights. It may be noted also that none of the correlations are as high as .6 while the means in all the tables are less than .4. Such coefficients are not considered high. The correlations corresponding for whole scales as given in Table C are +.71 for Rights and +.69 for Errors. Clearly then the cumulating of tests to

form what has been called a scale score has the effect of raising the correlation or, in other words, lengthening a test increases its reliability for these indexes. A more detailed discussion of this point will appear later.

Correlations between attempts on the nine components have been worked out for one group and are given in Table 28. Of the 36 coefficients, 27 are significant, the lowest average again occurring for Test 10 with each of the others. By three modes of indexing, then, this test shows up as distinct in type from the rest. The mean correlation for the whole table is $+.32$ which may be compared with the mean coefficient of $.56$ in Table 6. Lengthening the test also increases correlation when Attempts are employed as the index variable. The comparison is only a rough one, however, for somewhat different scales and groups are employed in the two cases.

TABLE 28.-- CORRELATIONS BETWEEN ATTEMPTS ON THE NINE OTIS COMPONENTS FOR I HIGH A

Test	1	2	4	5	6	7	8	9	10
1		$+.30$	$+.41$	$+.48$	$+.28$	$+.35$	$+.46$	$+.21$	$+.31$
2	$+.30$		$+.37$	$+.13$	$+.38$	$+.41$	$+.32$	$+.39$	$+.26$
4	$+.41$	$+.32$		$+.41$	$+.39$	$+.29$	$+.32$	$+.35$	$+.20$
5	$+.48$	$+.13$	$+.41$		$+.31$	$+.29$	$+.34$	$+.37$	$+.18$
6	$+.28$	$+.38$	$+.29$	$+.31$		$+.52$	$+.27$	$+.40$	$+.26$
7	$+.35$	$+.41$	$+.29$	$+.29$	$+.52$		$+.37$	$+.22$	$+.33$
8	$+.46$	$+.32$	$+.32$	$+.34$	$+.27$	$+.37$		$+.40$	$+.25$
9	$+.21$	$+.39$	$+.35$	$+.37$	$+.40$	$+.22$	$+.40$		$+.08$
10	$+.31$	$+.26$	$+.20$	$+.18$	$+.26$	$+.33$	$+.25$	$+.08$	
Mean	$+.34$	$+.31$	$+.32$	$+.31$	$+.34$	$+.35$	$+.34$	$+.29$	$+.24$

In Part I, Table 5, R and W showed a mean correlation of $-.49$ for whole scales. The corresponding coefficients for the Otis components are given in Table 60 with a mean of $-.21$. The increase in correlation by pooling is again evident. The correlations in the principal diagonal are between Rights and Errors on the same test and are therefore larger than the rest, the mean being $-.48$. The remainder of the table gives correlations for all possible combinations of Rights and Wrongs on the nine tests two at a time. All but three of the 31 coefficients are negative while just one third of them are significant according to the usual rule.

TABLE 60.- CORRELATIONS BETWEEN RIGHTS AND ERRORS ON THE NINE OTIS COMPONENTS FOR GRADE 7 AND I HIGH A

Wrongs	Rights									Mean
	1	2	3	4	5	6	7	8	9	
1	-.31	-.23	-.24	-.24	-.13	-.41	-.37	-.27	-.31	-.35
2	-.23	-.10	-.17	-.27	-.15	-.41	-.13	-.13	-.23	-.17
3	-.24	-.17	-.13	-.08	-.11	-.35	-.10	-.09	-.08	-.13
4	-.24	-.27	-.08	-.00	-.14	-.32	-.31	-.27	-.10	-.26
5	-.13	-.15	-.11	-.14	-.21	-.21	-.21	-.22	-.11	-.10
6	-.41	-.41	-.35	-.32	-.21	-.21	-.21	-.21	-.21	-.21
7	-.37	-.13	-.10	-.31	-.21	-.21	-.21	-.21	-.21	-.21
8	-.27	-.13	-.09	-.27	-.21	-.21	-.21	-.21	-.21	-.21
9	-.31	-.23	-.10	-.17	-.15	-.41	-.37	-.27	-.31	-.35
Mean	-.21	-.21	-.21	-.21	-.21	-.21	-.21	-.21	-.21	-.21

This suggests that the criterion of three times the probable error is too stringent for tests of this kind. If twice the probable error were adopted in the present case all coefficients over .13 would be significant including

five ninths of the total number, while all of those greater than one probable error or over .10 will include 56 out of 81. For the last case 25 coefficients are less than one probable error, yet 22 of them are negative in sign. Thus a coefficient less than one probable error appears to give assurance of negative correlation beyond the expectation from the usual rule of even chance; i.e. the probability of significance from the data appears to be greater than by theory. In any case, highly consistent negative correlation is exhibited by the whole array.

Section 15 Correlations of the Otis Components with Age

The correlations of the age factor with each of the Otis components are similar to those for whole scales. Table 30 shows higher correlations for the nine tests pooled than for the mean of the tests. Here again the effect of adding tests is to increase correlation. The formula for estimating the correlation of the pool of "n" tests with a criterion may be written in the form:

$$r_{(x_1+x_2+\dots+x_n)c} = \frac{r_{x_1c} + r_{x_2c} + \dots + r_{x_nc}}{\sqrt{n+2(r_{x_1x_2} + r_{x_1x_3} + \dots + r_{x_{n-1}x_n})}} \quad \begin{matrix} \text{(Theorem 4)} \\ \text{(Appendix)} \end{matrix}$$

Considering x_1, x_2, \dots, x_{10} as Rights for Otis Grade 7 and age as a criterion, the constants in Tables 24 and 30 give for this coefficient, $r_{(R_1+R_2+\dots+R_{10})age} = -.39$

a value identical with that obtained by pooling the nine components.

TABLE 30.- CORRELATIONS BETWEEN AGE AND THE INDEX VARIABLES ON THE OPIB SCALE FOR GRADE 7 AND I HIGH A

Test	Grade 7			I High A		
	Age x A	Age x R	Age x W	Age x A	Age x R	Age x W
1	-.13	-.32	+.13	-.32	-.33	+.22
2	-.29	-.41	+.18	-.27	+.31	+.26
3	-.32	-.34	+.27	-.19	-.29	+.21
4	-.33	-.42	+.33	-.33	-.47	+.25
5	-.33	-.33	+.33	-.31	-.52	+.33
6	-.33	-.33	+.33	-.31	-.36	+.19
7	-.33	-.33	+.33	-.13	-.41	+.33
8	-.33	-.33	+.33	-.33	-.44	+.24
9	-.33	-.33	+.33	-.33	+.21	-.64
Mean	-.35	-.36	+.25	-.29	-.37	+.27
All	-.33	-.33	+.27	-.33	-.33	+.34

Similar coefficients are given in Table 31. The differences between predicted and actual values are in no case significant. The above formula, then, appears to be a useful one in predicting the validity (correlation with a criterion) of tests by pooling components. It is to be noted also that the formula will have high values for large values of R_{xc} and small values of R_{xx} .

TABLE 31.- PREDICTED AND ACTUAL CORRELATIONS BETWEEN AGE AND INDEX VARIABLES

Group	Variables	Predicted Value	Actual Value
Grade 7	Rights Age	-.39	-.33
I High A	Rights Age	-.41	-.52
Grade 7	Wrongs Age	+.27	+.27
I High A	Wrongs Age	+.37	+.34
I High A	Attempts Age	-.33	-.32

To obtain a scale of high validity, therefore, component tests should be selected which have high correlation with the criterion but low correlation among themselves. Thorndike² justified the use of tests with low inter-correlations on the ground that they are repetitive; i.e. measures of the same fact. The above formula, however, will give high validity because the inter-test correlations occur in the denominator, and low values will thus raise the value of the fraction. The basis of test selection, then, appears to be mathematical, rather than psychological. Fortunately, however, the two bases agree.

Section 16 The Application of Reliability Formulae to Component Tests

In the preceding section it was shown that pooling component tests has the effect of increasing the validity of the total scale; i.e. to the extent to which it correlates with a criterion. The pooling of tests will next be shown to have a similar effect upon the reliability of a scale; i.e. the correlation between two forms of the same test.

It will be recalled that two forms of the Terman Scale were given to Group I High C on successive days. Reliability coefficients have been calculated for each of the 16 component tests and for all combined to give the total score. The results are given in Table 52.

TABLE 31.- RELIABILITY COEFFICIENTS FOR THE TER AN SCALE BY COMPONENTS AND TOTAL SCORE

Test	Formula	Correlation between Forms A and B	Rank
1	R	+.600	7
2	CR	+.870	3
3	R-W	+.300	6
4	R	+.900	1
5	CR	+.850	2
6	R-W	+.400	10
7	R	+.600	5
8	R-W	+.530	8
9	R	+.514	9
10	CR	+.700	4
Mean		+.670	
All		+.516	

Certain of the individual tests reveal a high degree of reliability, especially test 4 (logical selection) with a coefficient of .900. It will also be observed that some of the lowest correlations occur with tests scored R-W. This point will be dealt with more fully in a following section. The mean of the reliability coefficients on the tests is +.670, while the correlation for total Score on the two forms is +.516, so that pooling the tests has the effect of increasing reliability.

A predictive formula given by Brown,^a and also implied^b in Spearman's General Theorem may be given in the form,

$$R_{NN} = \frac{N R_{ii}}{(1+(N-1) R_{ii})} \quad (\text{Theorem 5, Appendix})$$

where R_{ii} is the correlation between two tests or the aver-

a. William Brown, *Essentials of Mental Measurement*, Cambridge University Press, London, 1911

b. C. Spearman, *Correlation of Sums and Differences*, Brit. Jour. of Psy., Vol. 5, 410-426

age of several, and N the number of tests thus amalgamated. In the present example the average correlation from the first three tests is .71. In order to predict the reliability coefficient for 10 such tests, it is only necessary to substitute these values in the above formula giving,

$$R_{(10)(10)} = +.96$$

The value from actual amalgamation is .92. Similarly, a calculation based upon the average of all 17 tests also gives

$$R_{(10)(10)} = +.96$$

The use of the formula in these cases, then, gives considerable over prediction.

In order to test the applicability of Brown's Formula more fully and to analyze more fully the effect of pooling tests on reliability, a more detailed procedure is next employed. Reliability coefficients on cumulated tests are obtained in two ways: The scores for tests 1 and 2 on each form of the Terman Scale are added, and the correlation determined; next tests 1, 2, and 3 are pooled and the two forms correlated, and so on until all 10 tests have been cumulated in this fashion. The second procedure is to begin with tests 10 and 2 and amalgamate in the reverse direction. These empirical results are then compared with theoretical values obtained by substituting $r = .68$ and N from 1 to 10 in Brown's Formula. Table 33 gives the results of this lengthy calculation.

TABLE 83.- THEORETICAL AND ACTUAL RELIABILITY COEFFICIENTS OBTAINED FROM BROWN'S FORMULA AND BY SUCCESSIVE CUMULATION OF THE TEST TRIAL COEFFICIENTS

Number of Tests Cumulated	Theoretical Value	Order of Cumulation	
		1 to 10	10 to 1
1	+.63	+.61	+.70
2	+.57	+.61	+.72
3	+.47	+.57	+.68
4	+.39	+.51	+.66
5	+.33	+.42	+.61
6	+.26	+.32	+.56
7	+.21	+.23	+.51
8	+.16	+.14	+.47
9	+.10	+.08	+.37
10	+.06	+.02	+.22

Inspection of the table shows a rough agreement in the three series. Figure 7 which is based on the table brings out the comparisons more clearly. The three curves show a rapid initial rise up to 4 cumulated tests and then a more gradual increase to the various values. The more rapid rise of the curve cumulated from tests 1 to 10 than the reverse one, is no doubt due to the greater reliability of the first few tests as indicated in Table 82.

While the general agreement in the three curves is evident, nevertheless there is a very clear tendency for the theoretical curve calculated for $r = .31$ to give an over-prediction beyond 4 or 5 cumulated tests. This may be partly due to the unequal values of the individual reliability coefficients given in Table 82, but whatever the cause, the use of Brown's Formula for prediction in a case of this kind is open to question. The equation

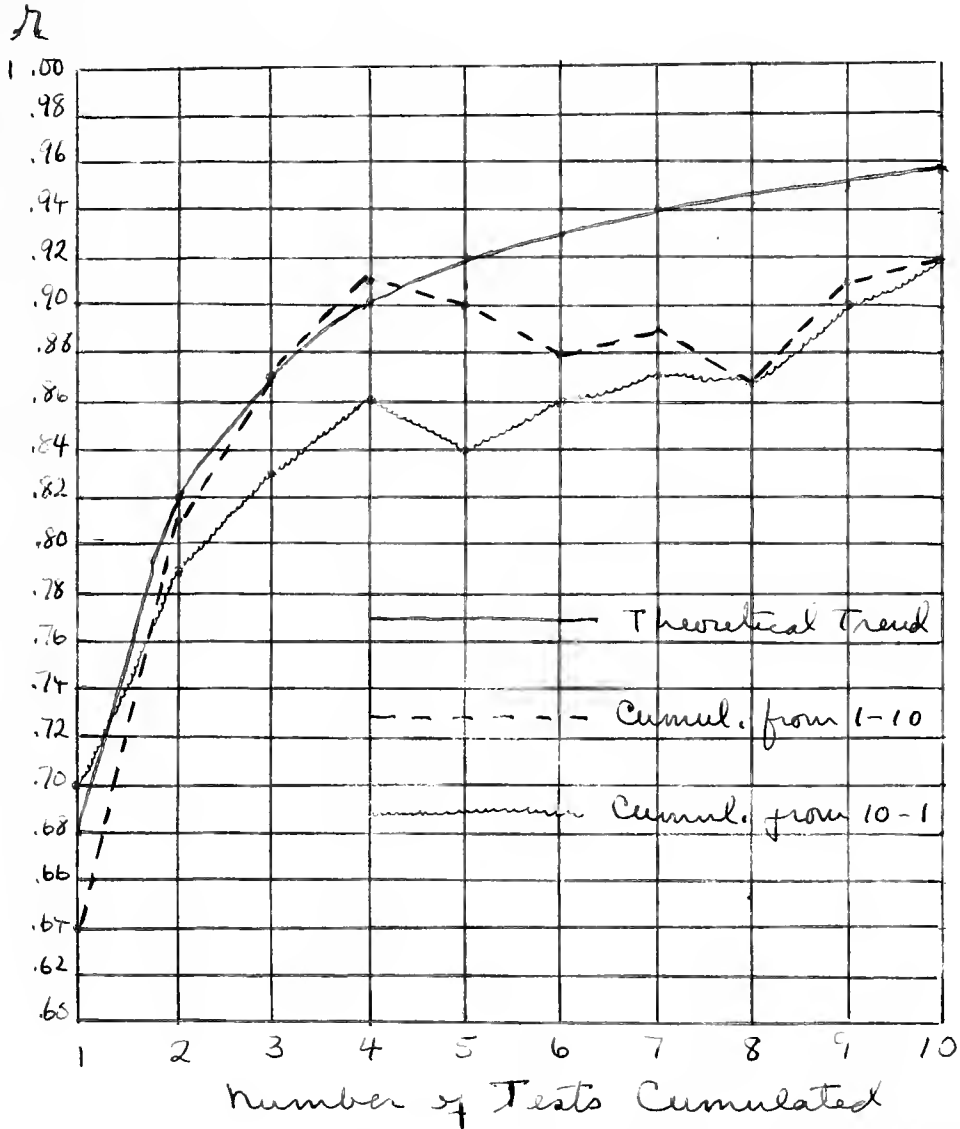


Fig 5. Theoretical and Actual Reliability Trends based on Tin Terman Component

and corresponding theoretical curve indicate that to get any desired degree of reliability with $+1.00$ as an upper limit, it is only necessary to amalgamate tests indefinitely. This is, of course, absurd. The formula gives an over prediction fairly early in the series of cumulated tests. From the above tables it appears that four or five typical tests of the battery will give almost as reliable an index as the pool of all ten components. This result would account in part for the high reliability of such tests as the Chicago Scale consisting of only five components.

This problem is one of great importance in test construction. If intelligence can be indexed with almost as great accuracy by a short scale as by one twice as long, the saving in time alone is enormous. Moreover if the short series can be shown to be as valid as the longer one by correlation with criteria, the abbreviated method is further justified. Inasmuch as no suitable criterion other than age was available for the present data, the check cannot be rigidly applied. The age correlations by half and by whole scales, however, agree almost exactly ($-.37$, $-.39$), so that with age as a criterion, the five test battery is as valuable as the ten test scale.

Section 17 Summary of Analysis of Components

The relationships found between index variables in Part I are verified for component tests. These coefficients are in general lower than by whole tests, so that pooling has the effect of increasing the correlation between indexes. Inter-correlations between components for R, W, and A reveal a high degree of consistency for such short tests but are less stable than for similar coefficients by whole tests. Furthermore the consistent batteries of correlations even for R and W on different tests indicate a high degree of homogeneity in the test material with the possible exception of Test 10.

In addition to raising the correlation between index variables, pooling tests also has the general effect of increasing the validity and reliability of a scale within certain limits. Predictive formulae are useful in this connection but are likely to give an over-estimate of the correlation to be expected by pooling. Moreover the physical endurance of the children determines the maximum length of the tests at a sitting, so that the formulae are limited in application. The gain in validity and reliability is rapid on pooling the first few tests, but the point is soon reached where the addition of similar material affects the correlations but slightly. The results indicate that a battery of four or five carefully selected components will give an index with substantially the same reliability as a scale twice that length.

Part III Scoring FormulaeSection 18 The Linear Form, $S = a (R + bW)$ a. Formulae with Highest Validity

In Parts I and II it has been shown that the scoring formulae employed by the authors of the scales have little effect upon the resultant scores when a number of components are pooled. The Terman Scale with the components scored by the three formulae, $S = R$, $S = 2R$, and $S = R + W$ has a correlation of +.26 with the score obtained by using $S = R$ on all ten components (Tables 4 and 5). The amalgamated score then, is not very sensitive to such changes in the component scoring formulae and simple forms are recommended on these grounds. The single component, however, is much more violently affected by changes in the formulae employed to index it. Changes in weights which affect the pooled score but slightly, will be found to have a pronounced effect upon the individual components.

Table 4 which gives the various component scoring formulae used by the authors of the scales, includes only formulae of the linear type; i.e. equations of the first degree in the variables employed. These variables are R and W , so that the most general formula used may be written,

$$(1) \quad S = a(R + bW) = aR + cW$$

where a , b , and c , are constants. It was also noted in Section 3 that the relationship, $A = R + W$

makes it possible to express this formula in terms of R and A or W and A. Formula (1), however, has been so generally employed that it will be adopted here for further analysis. Formulae expressed in terms of the other variables may be obtained by substitution if they are required.

The question immediately arises as to the best values to assign the constants a and c in equation 1. A general solution of this problem may be obtained by the method of least squares. Values for R and W are obtained for each of the N individuals of a given population. Assuming that a criterion, K, is the best measure of such determination, a set of N equations may be formed,

$$\begin{aligned} K_1 &= a_1 R_1 + c_1 W_1 \\ K_2 &= a_2 R_2 + c_2 W_2 \\ &\text{---} \text{---} \text{---} \text{---} \text{---} \text{---} \\ K_N &= a_N R_N + c_N W_N \end{aligned}$$

where the K's, R's, and W's are known, and a's and c's are to be determined so as to minimize the inconsistency in the equations which is assumed to be due to imperfect measurement.

Next V_1, V_2, \dots, V_N will be written for the differences between K_1, K_2, \dots, K_N and the values obtained from the best determinations for the a's and c's:

$$\begin{aligned} a_1 R_1 + c_1 W_1 - K_1 &= V_1 \\ a_2 R_2 + c_2 W_2 - K_2 &= V_2 \\ &\text{---} \text{---} \text{---} \text{---} \text{---} \text{---} \\ a_N R_N + c_N W_N - K_N &= V_N \end{aligned}$$

These differences or "residuals" are assumed to be

normally distributed. While the assumption is open to question for data of this type, it is nevertheless the best that can be made. The most probable values for a and c next require that the sum

$$V_1^2 + V_2^2 + \dots + V_N^2 = \text{a minimum.}$$

The remainder of the procedure consists in setting up the "normal equations" in the usual way. Transferring the variables to their respective means, and setting up these equations gives,

$$a \sum R^2 + c \sum RW = \sum KR$$

$$a \sum RW + c \sum W^2 = \sum KW$$

Since, $\sigma_x^2 = \frac{\sum X^2}{N}$, and $\rho_{xy} = \frac{\sum xy}{N \sigma_x \sigma_y}$, these equations may be written in the form,

$$a \sigma_R + c \rho_{RW} \sigma_W = \rho_{KR} \sigma_K$$

$$a \rho_{RW} \sigma_R + c \sigma_W = \rho_{KW} \sigma_K$$

Solving these equations for a and c gives,

$$(3) \quad a = \frac{\sigma_K (\rho_{KW} \rho_{RW} - \rho_{KR})}{\sigma_R (\rho_{RW}^2 - 1)}$$

$$(4) \quad c = \frac{\sigma_K (\rho_{KR} \rho_{RW} - \rho_{KW})}{\sigma_W (\rho_{RW}^2 - 1)}$$

The value $\frac{c}{a}$ may then be written

$$(5) \quad C = \frac{c}{a} = \frac{\sigma_R (\rho_{KR} \rho_{RW} - \rho_{KW})}{\sigma_W (\rho_{KW} \rho_{RW} - \rho_{KR})}$$

This last result has been obtained by Thurstone² as a

a. L.L. Thurstone, A Scoring Method for Mental Tests, Psy. Bull., Vol. XVI, No.7, July, 1919.

value for C in the formula $S = R + CW$ such that the correlation ρ_{ks} is a maximum; i.e. C or c is determined in such a way as to possess the highest validity with a criterion. The formula for this correlation is,

$$(6) \quad \rho_{ks} = \frac{\rho_{kn}\sigma_n + C \rho_{kw}\sigma_w}{\sqrt{\sigma_n^2 + 2C \rho_{nw}\sigma_n\sigma_w + C^2\sigma_w^2}}$$

Thurstone also makes use of Yule's equation for multiple correlation to obtain an expression for the highest correlation with the linear formula $S = a + CW$. This result may be written

$$(7) \quad R_{k(a+cw)} = \sqrt{\frac{\rho_{kn}^2 + \rho_{kw}^2 - 2 \rho_{kw} \rho_{kn} \rho_{nw}}{1 - \rho_{nw}^2}}$$

The actual procedure involved in determining the constants a and c will then be as follows:

1. Give the test to a group and obtain also the criterion of validity against which the formula is to be checked.
2. Score the test for R and W and compute the constants ρ_{kn} , ρ_{kw} , ρ_{nw} , σ_n , and σ_w (σ_k is not required if Formula (6) is employed).
3. Substitute these last results in equations (6), (4), and (7) and obtain the formulae $S = aR + cW$ or $S = R + \frac{c}{a}W$ (differing only by factor of no arbitrariness, $\frac{1}{a}$).
4. To predict the highest correlation obtainable with these formulae, substitute the computed constants in equation (7).

b. Limitations in the Use of the Formula $S = aR + cW$

In section 3, two plans for administering tests were described. According to the first, the time is fixed and Attempts, Rights, and Wrongs allowed to vary; according to the second, the number of Attempts is fixed, while Time, Rights, and Wrongs are recorded. Two of the index variables are thus alternately controlled by the method of administering the test.

The formula $S = aR + cW$, with constants determined by highest validity with the criterion, serves very well for tests given according to the plan of fixing the time. All of the tests employed in the three intelligence scales are of this type and hence no difficulty is encountered.

For tests administered by fixing the Attempts, however, the above formula is inadequate. This arises from the fact that S is no longer a function of two independent variables, but of only one. If Attempts are constant and $A = R + W$, then $W = -R + d$ where d is ^aconstant. Substituting the value for W in the expression $S = aR + cW$ gives $S = aR + c(d - R)$ or $S = (a - c)R + cd$. Thus the scoring formula is independent of W , and no matter what value is assigned to c (with the exception of $c = a$, for which value the correlation is zero) the correlation r_{KS} given by formula (6) is equal to r_{KR} . This last result follows from the fact that $r_{x(ay+b)} = r_{xy}$ where a and b are constants (Theorem 1, Appendix)

In other words if Attempts are constant, the scoring formula $S = a R$ has the same validity with a criterion as any linear function of R and W (with the exception $R + W$ for which $Aks = 0$).

A further limitation in the use of the formula $S = R + CW$ lies in its sensitiveness. The value of C as determined by Formula (5) depends upon σ_R and σ_W which in turn depend upon the proportion of Rights and Wrongs in the group. Thus while the value for C is the best prediction for the particular group, the value will very likely differ very materially in other groups with differing percentages of Wrongs. This means that the formula $S = R + C \cdot W$ may be used with safety only (for tests administered with time fixed) when the value for C has been determined after the test has been given to the particular group, not before. Thurstone^a discusses the sensitiveness of the formula, but is unwilling to admit its limitations. Table 34 which is based on extremely lengthy and careful calculation, indicates wide variation in the determinations of C for the 9 Otis components on Grade 7 and I High A. The average difference in the values of C for the two groups is .649 and in only one test is the difference less than 1. These results point clearly to the conclusion that the above formula has little general merit; i.e. the value for C from one group cannot be safely assumed to hold for another.

TABLE 34.- VALUES FOR C IN THE FORMULA $S = R + CW$ DETERMINED FOR GROUPS GRADE 7 AND I HIGH A ON THE NINE COMPONENTS OF THE OPTI SCALE (AGE GRADATION)

Test	Grade 7						I High A		
	$Age \times R$	$R \times W$	$Age \times W$	\overline{C}_R	\overline{C}_W	C_1	$Age \times R$	$R \times W$	$Age \times W$
1	-.381	-.823	+.126	2.631	2.520	+.317	-.327	-.767	+.017
2	-.482	-.170	+.182	2.812	1.794	-.414	-.300	-.047	+.282
4	-.277	-.122	-.066	2.122	2.232	+.001	-.220	-.752	+.312
5	-.417	-.662	+.342	2.477	2.429	+.109	-.402	-.627	+.230
6	-.370	-.430	+.042	2.742	2.470	+.072	-.302	-.652	+.375
7	-.261	-.276	+.214	2.738	2.052	-.654	-.256	-.518	+.189
8	-.332	-.070	+.072	2.852	2.027	+.125	-.402	-.837	+.327
9	-.122	-.274	+.000	2.110	2.507	+.370	-.420	-.412	+.244
10	+.080	-.040	-.347	2.222	2.372	2.270	+.000	-.051	-.040

TABLE 34- CONTINUED

Test	I High A			Diff.
	\overline{C}_R	\overline{C}_W	C_2	
1	1.212	0.100	+.035	-.355
2	0.034	0.000	+.182	+.875
4	1.165	0.005	+.000	+.000
5	2.572	2.427	+.0272	-.112
6	2.334	2.520	+.000	+.000
7	2.271	2.372	+.0207	-.047
8	2.330	2.638	+.0100	+.000
9	2.328	2.710	+.0372	+.040
10	2.256	2.372	+.036	+.140

c. The Use of the Formula $S = R + CW$

In test material where two alternatives are given for each item and guessing therefore possible the formula

$S = R - W$ (i.e. $C = -1$) has been frequently adopted. Similar formulae are used when the number of choices is greater. These expressions are assumed to correct for guessing element involved. According to the above equation a person guessing blindly on all of the items will get half of them right by chance, and hence a zero score which he deserves.

For actual guessing, then, the form -

p. 67

mula penalizes justly; but it also penalizes for errors which are not due to guessing, and hence unjustly. As a result, for very difficult material, nearly all of the scores may be negative.

A number of experimental attempts have been made to determine the amount of guessing in tests of this sort. After administering a set of True-False tests of the syllogistic reasoning type, the writer asked the pupils on which items they had guessed. The number guessed was about two per cent of the total number of errors made, and of those items guessed only 30 per cent were wrong. For such groups and tests, the penalty attached to errors by the formula $R-W$ is enormously too great. It may be noted, however, that the children did not always know whether or not they had guessed on any item. Reasoning and guessing are often indistinguishable, and who has not credited himself with reasoning when he has only made a lucky guess.

Instead of assuming that a penalty should be attached for guessing, Thurstone proposes to use the formula with the value of C to be determined according to validity as above. This method appears to be preferable to that of a priori determination, when tests are administered with the time fixed. If Attempts are fixed, however, the formula becomes independent of W as has just been shown, and all values of C (except $C=+1$) gives the same correlation

$$R_{K(R+CW)} = R_{KR} \quad .$$

This point is of great practical importance because most tests of the True-False type are administered so that all may finish; i.e. Attempts constant. According to the above results, all formulae of the type $S = R - \frac{1}{n} W$ give the same correlation with a criterion as is obtained by using Rights alone; i.e. $S = R$. Inasmuch as the expression $R - \frac{1}{n} W$ does not correct adequately for guessing, and has the same validity as $S = R$ for Attempts constant, the writer believes it should be abandoned in favor of the simpler form.

The following example illustrates the foregoing discussion.

K	R	W	R-W	K	n	w	n-w
10	3	8	-5	-30	-3	3	-4
20	3	8	-5	-10	-3	3	-4
30	5	5	0	0	1	-1	2
40	4	6	-2	10	0	0	0
50	7	3	4	20	3	-3	6
<u>150</u>	<u>20</u>	<u>30</u>	<u>-10</u>				
50	4	6	-2				

K^2	n^2	$(n-w)^2$	Kn	$K(n-w)$
400	4	16	40	80
100	4	16	20	40
0	1	4
100	1	0
400	9	36	60	120
<u>1000</u>	<u>18</u>	<u>72</u>	<u>120</u>	<u>240</u>

$$r_{kr} = \frac{120}{\sqrt{18000}} = \frac{2}{\sqrt{5}}$$

$$r_{k(n-w)} = \frac{240}{\sqrt{72000}} = \frac{2}{\sqrt{5}}$$

Section 1.2. Simple Ratios

a. The Correlation Between Speed and Accuracy

The formulae $S = \frac{A}{T}$ and $S = \frac{R}{A}$ may be conveniently employed to index Speed and Accuracy, respectively. The latter form has been used extensively in the first two parts of this study and has been found to have valuable properties which other indexes do not possess. With Time constant as in the intelligence series discussed, the formula for Speed reduces to $S = \frac{1}{T}$; i.e. the Attempts give the measure of Speed directly when Time is fixed. When Attempts are constant, the Speed is given by the reciprocal of the Time (in suitable units).

The general expression for the correlation between two ratios $\frac{X}{Y}$ and $\frac{Z}{W}$ may be written in the form

$$(6) \quad \rho_{\frac{X}{Y}, \frac{Z}{W}} = \frac{\rho_{XZ} V_X V_Z + \rho_{YW} V_Y V_W - \rho_{XW} V_X V_W - \rho_{YZ} V_Y V_Z}{\sqrt{[V_X^2 - 2\rho_{XY} V_X V_Y + V_Y^2]} [V_Z^2 - 2\rho_{ZW} V_Z V_W + V_W^2]}$$

where the V 's are coefficients of variation given by the formula $V = \frac{100 S.D.}{M}$ (Theorem 2 Appendix). For the two ratios $\frac{A}{T}$ and $\frac{R}{A}$ this expression becomes

$$\rho_{\frac{A}{T}, \frac{R}{A}} = \frac{\rho_{AR} V_A V_R + \rho_{TA} V_T V_A - \rho_{AA} V_A V_A - \rho_{RT} V_R V_T}{\sqrt{[V_A^2 - 2\rho_{AT} V_A V_T + V_T^2]} [V_R^2 - 2\rho_{RA} V_R V_A + V_A^2]}$$

For $T = \text{constant}$, $V_T = \text{constant}$, and all correlations with T are zero, therefore,

$$\rho_{\frac{A}{T}, \frac{R}{A}} = \frac{\rho_{AR} V_A V_R - V_A^2}{\sqrt{V_A^2} [V_R^2 - 2\rho_{RA} V_R V_A + V_A^2]}$$

which expression reduces to

$$(2) \quad \rho_{\frac{A}{T}, \frac{R}{A}} = \frac{\rho_{AR} V_R - V_A}{\sqrt{V_R^2 - 2\rho_{RA} V_R V_A + V_A^2}}$$

Equation (9) thus gives the correlation between Speed and Accuracy for tests where Time is fixed. The Maximum value, or $r_{AR} = +1.00$ is given for $r_{AR} = +1.00$. Thus if the pupils get every problem that they attempt right, Speed and Accuracy will be perfectly correlated. For zero correlations between Attempts and Rights, Speed and Accuracy are negatively correlated, the value approximating $-\frac{\sqrt{2}}{2}$. Finally, if the ratio $\frac{V_A}{V_R}$ is equal to the value of r_{AR} , the correlation between Speed and Accuracy will be zero. Tables 5 and 12 indicate that these last relationships will hold very closely for intelligence test data. The ratios of the V_A' and the corresponding correlations are approximately .7

Formula (9) is very useful for obtaining the correlation between Speed and Accuracy from the single correlation table for Attempts and Rights. This table gives the values $r_{RA}, r_{R,GA}, r_{MR},$ and r_{MA} which are all that are required to obtain r_{AR} . The correlations with Accuracy in Part I were computed from ratios $\frac{R}{A}$ obtained for each variate by division. About a third of the coefficients were then checked by the above formula. Substitution in equation (9) requires but a few moments, while the division for ratios alone takes about an hour for 50 cases. A great saving of time is, therefore, effected by the use of the above formula, especially if the constants in the formula are needed for other purposes.

If Attempts are fixed, Speed is measured by the recip-

rocal of the Time, and Accuracy by Rights. Equation (8) then, reduces to

$$r_{\frac{A}{T} \frac{R}{A}} = -r_{RT}$$

For reasoning test material of the True-False type, administered with Attempts constant, low negative correlations were obtained for Rights and Time, indicating that the correlation between Speed and Accuracy given by the last formula is positive and low. For 15 groups of about 20 pupils each, the average correlation $r_{\frac{A}{T} \frac{R}{A}} = .20 \pm .05$. For both types of test administration, then, Speed and Accuracy exhibit correlation that is zero or barely large enough to be significant.

b. The Validity of Simple Ratios as Scoring Indexes

It has just been shown that the ratios giving Speed and Accuracy are relatively independent measures of intelligence. The choice of the proper index will therefor depend upon criteria such as the purpose for which the measurements were made, the validity of the index, and its reliability. The question of validity will be taken up first.

The use of certain linear formulae has been justified on the basis of their validity or correlation with a criterion. This same principle may be applied to ratios. If a criterion, K , be substituted in Formula (8) in place of $\frac{Z}{W}$ an expression for the validity of $\frac{X}{Y}$ may be written in the form

$$(10) \quad r_{K \frac{X}{Y}} = \frac{r_{KX} V_X - r_{KY} V_Y}{\sqrt{V_X^2 - 2r_{XY} V_X V_Y + V_Y^2}}$$

The correlation tables for R_{KX} and R_{KY} will furnish all of the data necessary for this formula and save the labor of calculating $R_{K\frac{X}{Y}}$ by the direct method of division. In general if $R_{K\frac{X}{Y}}$ is significantly higher than R_{KX} or R_{KY} its use in place of the single variables is justified on the grounds of higher validity.

While score is not the best measure of validity for indexes which are so closely related to it, nevertheless the data in Table 5 will give a suggestion as to the general method. The order of the correlations between $R, \frac{R}{A}, A$, and score is $R_{SR} = .96$, $R_{S\frac{R}{A}} = .76$, $R_{SA} = .59$. The differences are all significant here so that on the basis of validity $\frac{R}{A}$ is a better index than A , but not so good as R alone.

c. The Reliability of Bin in Tables of Scoring Indexes

Equation (3) is again useful in predicting the reliability of a ratio without direct calculation on the ratios themselves. If $\frac{X_1}{Y_1}$ and $\frac{X_2}{Y_2}$ denote the ratios in question on successive trials of the same test or b, parallel forms, the reliability formula may be written in the form

$$(11) \quad R_{\frac{X_1}{Y_1}, \frac{X_2}{Y_2}} = \frac{R_{X_1X_2}V_{Y_1}V_{Y_2} + R_{Y_1Y_2}V_{X_1}V_{X_2} - R_{X_1Y_2}V_{X_1}V_{Y_2} - R_{X_2Y_1}V_{X_2}V_{Y_1}}{\sqrt{[V_{X_1}^2 - 2R_{X_1Y_1}V_{X_1}V_{Y_1} + V_{Y_1}^2][V_{X_2}^2 - 2R_{X_2Y_2}V_{X_2}V_{Y_2} + V_{Y_2}^2]}}$$

In order to calculate this quantity, four correlation tables are required: X_1X_2 , X_1Y_2 , X_2Y_1 , and X_2Y_2 .

If they are prepared all at once, the marginal frequencies give excellent checks on the distributions. As in the case of validity if $R_{X_1 X_2}$ is significantly greater than $R_{X_1 X_2} \sim R_{X_2 X_1}$, it is to be preferred as the more reliable index. This method was applied to the Ternan Scale Forms A and B. The results of the direct calculation appear in Table 7. Prediction by Formula (11) gives: $R_{R_1 R_2} = +.896$, $R_{A_1 A_2} = +.854$, and $R_{A_1 R_2} = +.684$. Accuracy and Rights are significantly more reliable than Attempts, but are not essentially different from one another.

The intelligence quotient, and similar ratios, are essentially a score divided by a chronological age, when the score is expressed in age units and taken from a suitable origin. The choice of the age unit is important chiefly because the resulting ratio is then a pure number and easily interpreted. As far as the validity and reliability of the ratio are concerned, the choice of the unit for score is of no consequence, inasmuch as correlation is a measure independent of the units employed for the two variates. The origin from which the score is taken is always of importance since any shift obviously changes the ratio e.g. $\frac{X}{A} \neq \frac{X+C}{A}$. The origin for the score in the intelligence quotient is the same as for chronological age, hence no difficulty is encountered.

Formula (8) is a function of R 's and \sqrt{C} , only the lat-

ter being affected by the origin from which the variables are taken. Indeed, by suitable choice of origin all of the V_s may be made equal, so that the formula reduces to

$$(12) \quad \rho_{\frac{X'Z'}{Y'W'}} = \frac{\rho_{X'Z'} + \rho_{Y'W'} - \rho_{X'W'} - \rho_{Y'Z'}}{2\sqrt{(1 - \rho_{X'Y'})(1 - \rho_{Z'W'})}}$$

Letting Y' and W' denote age, and X' and Z' scores on successive trials of a test, an expression of the reliability of the ratio $\frac{S}{age}$ may be written in the form

$$\rho_{\frac{S_1}{age} \frac{S_2}{age}} = \frac{1 + \rho_{S_1 S_2} - \rho_{S_1 age} - \rho_{S_2 age}}{2\sqrt{(1 - \rho_{S_1 age})(1 - \rho_{S_2 age})}}$$

Furthermore if $\rho_{S_1 age} = \rho_{S_2 age}$ this expression reduces to

$$(13) \quad \rho_{\frac{S_1}{age} \frac{S_2}{age}} = \frac{1 + \rho_{S_1 S_2} - 2\rho_{S age}}{2 - 2\rho_{S age}}$$

Formula (13) will have the value +1.00 for $\rho_{SS} = +1.00$, and will be zero for $1 + \rho_{S_1 S_2} = 2\rho_{S age}$. For a given positive value of $\rho_{S_1 S_2}$, the reliability given by the formula will increase as $\rho_{S age}$ decreases from the value $\frac{1 + \rho_{S_1 S_2}}{2}$.

As an illustration, let $\rho_{S_1 S_2} = +.8$. The function $\rho_{\frac{S_1}{age} \frac{S_2}{age}}$ may then be tabled for the argument $\rho_{S age}$ as follows:

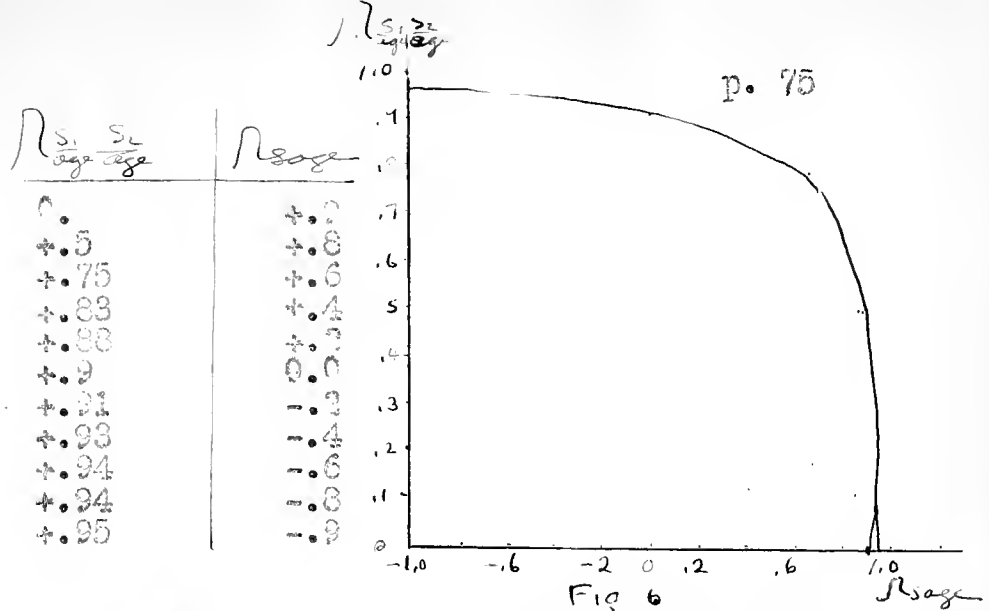


Figure 6 which is based on the table shows that the highest reliability of the ratio $\frac{S}{age}$ occurs with the highest negative values of r_{sage} and decreases as r_{sage} increases (to the right). The value for which $r_{\frac{S_1}{age} \frac{S_2}{age}} = .8$ is found by solving the equation

$$.8 = \frac{1.8 - 2x}{2 - 2x}$$

giving $x = .5$, so that the ratio has a greater reliability than $r_{S_1 S_2}$ up to the value $r_{sage} = .5$. Inasmuch as .8 is a good reliability coefficient for score, and r_{sage} is seldom as high as +.5, the supposititious example above indicates that ratios with age have in general a greater reliability than that between crude scores. The negative correlation usually found between score and age for a given grade group (See Tables 8 and 9) also indicates that ratios of the above type are most reliable when the group is thus selected with respect to age.

A final example will be given to illustrate the effect of transferring the origin to eliminate the . Assuming the approximate values $VA = VB = 20$, $r_{AB} = .9$, $r_{age} = .8$ and $r_{sage} = -.4$, Formula (11) will give by simple substitution

$$r_{\frac{A}{age} \frac{B}{age}} = +.93$$

Shifting the origin to eliminate the V 's and applying formula (13) gives

$$r_{\frac{A}{age} \frac{B}{age}} = +.96$$

p. 76

An increase in reliability of .73 is thus brought about by the transfer of origin. For the linear form $S = aR + cW$ such a shift will, of course, have no effect upon the reliability or validity of the formula (See Theorem 1 Appendix).

The bases for selecting a simple ratio as a mode of indexing may be listed as follows:

1. The desirability of indexing another feature of the general characteristic, e.g. Accuracy, though a sub-characteristic of intelligence is essentially different from Speed.
2. The general properties of the ratio as compared with other variables e.g. the intelligence quotient facilitates comparisons with normal achievements.
3. The validity of the ratio to be determined by Formula (12).
4. The reliability of the ratio as determined by formula (11).

Section 2^o General Conclusions

1. The various types of response to test material have been treated as index variables for the traits in question. An analysis of these variables for intelligence test data

revealed fairly definite relationships between them as indicated by the coefficient of correlation.

2. By eliminating the difficulty factor, the primary index variables were reduced to A, R, W, and T, one of these being fixed by the method of test administration.

3. An analysis of whole scales indicated that all of the primary variables have valuable properties as indexes. The introduction of the simple ratio made possible a comparison of the indexes revealing them in order of general reliability as S, R, R, W, and A, with Time fixed.

4. According to the criteria employed, the complicated formulae used by the authors of the tests, are not significantly better than Rights alone. Accuracy comes next as a generally reliable index, and has properties which the other variables do not possess. For batteries of tests, then, scoring by Rights alone is justified by reason of greater simplicity and practically equal reliability as compared with more complicated formulae.

5. In discriminative capacity, Attempts proved to be highest and Accuracy lowest; i.e. individuals and groups differ more widely in Speed than in Accuracy. Lack of discrimination between groups is of less consequence in an index than failure to differentiate between individuals. Accuracy, therefore, retains its high place as an

index regardless of the slight inter-group differences shown.

6. Analysis of the scales by component tests furnished a check upon the results obtained for whole batteries. The coefficients in general are lower and less consistent than by whole scales.

7. The validity and reliability of tests are increased by pooling components. Estimations of these correlations are furnished by certain predictive formulae, which in general, tend to give over estimation very early in the cumulated series.

8. Both theoretical prediction and actual results indicate that pooling tests soon ceases to increase validity and reliability materially. A battery of five well-selected tests is about as satisfactory as one twice as long.

9. The formula $S = aR + cW$ has been shown to be the most general linear form. Formulae for validity have been worked out by the method of least squares.

10. The linear formula above is open to ^{objection}~~question~~ because of its sensitiveness for values of C . This implies that a new determination is necessary for each new group dealt with.

11. The scoring formula $R-W$ is criticized because of its

failure to correct for guessing. It has been shown that for Attempts constant (the usual method with True-False tests) corrective formulae of the type $R - \frac{1}{K} W$ have exactly the same validity as Rights alone.

12. Special formulae and methods have been worked out for determining the validity and reliability of simple ratios.

13. The valuable properties possessed by simple ratios indicate that they are highly desirable and useful scoring devices in spite of the labor of division. Special tables give such quotients directly.

14. More complicated formulae have not been dealt with because the labor involved in their use would be prohibitive¹ no matter what virtues they might be found to possess. The sensibleness of such formulae is sufficient reason for avoiding them.

15. The results as a whole point to the conclusion that for batteries, and for simple tests as well, the most desirable and useful indexes are R and $\frac{R}{K}$.

Bibliography

Memoirs of the National Academy of Sciences, Vol. XV,
Psychological Testing in the United States Army.

C.Spearman, Correlation of Sums and Differences, British
Journal of Psychology, Vol. V, pp. 412-426.

A.S.Otis, An Absolute Point Scale for the Group Measure-
ment of Intelligence. Journal of Educational Psychology,
Vol. II, Nos.5 and 6, May-June, 1918.

L.L.Thurstone, A Scoring Method for Mental Tests, Psy.
Bull., Vol.XVI, No. 7, July 1919.

T.L.Kelley, The Reliability of Test Scores JI.Ed.Research
Vol.III, No.5, p.379, May 1931

A.S.Otis and H.E.Knollin, The Reliability of the Binet
Scale and Pedagogical Scales, JI.Ed. Research., Vol. IV, No.2
121, Sept. 1931

Bibliography --Texts

C.V. Rieu, An Introductory Course in Mathematical Statistics,
Charles Griffin and Company, London, 1911.

A.L. Bowley, Elements of Statistics, Charles Scribner's
Sons, New York, 1907

D.C. Jones, A First Course in Statistics, G. Bell and Sons
London, 1903

H. Brown and C.H. Thompson, The Essentials of Mental
Measurements, Cambridge University Press, London, 1924

APPENDIX A

Correlation Tables for Reliability Coefficients

CORRELATION TABLE FOR SCORES ON TRIAL FORM 1 WITH SCORE ON TRIAL FORM 3 WITH GROUP 1 HIGH C

SCENE ON TIERMAN FORT B (FIRST TRIAL)

[illegible]

$$\mu = +.908 \pm .010$$

$$M_A = 101.74 \pm 1.79$$

$$M_B = 84.11 \pm 1.96$$

$$M_A - M_B = 17.63 \pm 2.65$$

$$\sigma_B = 33.79 \pm 1.39$$

$$\overline{G}_A = 30.91 \pm 1.27$$

$$\sigma_B - \sigma_A = 2.88 \pm 1.88$$

∴ gain is significant

diff. may be significant

CORRELATION TABLE FOR ATTEMPTS ON TERMAN FORM A WITH
ATTEMPTS ON TERMAN FORM B WITH CORRELATION COEFFICIENT

ATTEMPTS ON TERMAN FORM B (FIRST GROUP)

ATTEMPTS ON TERMAN FORM A	80-89	90-99	100-109	110-119	120-129	130-139	140-149	150-159	160-169	170-179	180-189	f
180												1
170												1
160												1
150												1
140												1
130												1
120												1
110												1
100												1
90												1
80												1
70												1
60												1
50												1
40												1
30												1
20												1
10												1
0												1
f	3	3	4	10	11	10	17	21	20	20	8	135

$$r = +.684 \pm .032$$

$$\begin{aligned} M_A &= 161.00 \pm 1.20 \\ M_B &= 148.19 \pm 1.40 \\ \hline M_A - M_B &= 12.81 \pm 1.84 \end{aligned}$$

$$\begin{aligned} \sigma_B &= 24.12 \pm 0.99 \\ \sigma_A &= 20.63 \pm 0.85 \\ \hline \sigma_B - \sigma_A &= 3.49 \pm 1.30 \end{aligned}$$

∴ gain is significant

∴ diff. sig. is significant

SIGHTS ON TRIUMPH BOULEVARD (FIRST TRIAL)

$$\mu = +.896 \pm .011$$

$$\sigma_B - \sigma_A = 3.41 \pm 1.47$$

did not become significant

CORRELATION TABLE FOR WRONGS ON TERMAN FORM A WITH
WRONGS ON FORM B WITH GROUP I HIGH C

WRONGS ON TERMAN FORM B (FIRST TRIAL)

WRONGS ON TERMAN FORM A											
	10-19	20-29	30-39	40-49	50-59	60-69	70-79	80-89	90-99	100-109	110-119
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0
20	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0	0	0
24	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0	0	0
27	0	0	0	0	0	0	0	0	0	0	0
28	0	0	0	0	0	0	0	0	0	0	0
29	0	0	0	0	0	0	0	0	0	0	0
30	0	0	0	0	0	0	0	0	0	0	0
31	0	0	0	0	0	0	0	0	0	0	0
32	0	0	0	0	0	0	0	0	0	0	0
33	0	0	0	0	0	0	0	0	0	0	0
34	0	0	0	0	0	0	0	0	0	0	0
35	0	0	0	0	0	0	0	0	0	0	0
36	0	0	0	0	0	0	0	0	0	0	0
37	0	0	0	0	0	0	0	0	0	0	0
38	0	0	0	0	0	0	0	0	0	0	0
39	0	0	0	0	0	0	0	0	0	0	0
40	0	0	0	0	0	0	0	0	0	0	0
41	0	0	0	0	0	0	0	0	0	0	0
42	0	0	0	0	0	0	0	0	0	0	0
43	0	0	0	0	0	0	0	0	0	0	0
44	0	0	0	0	0	0	0	0	0	0	0
45	0	0	0	0	0	0	0	0	0	0	0
46	0	0	0	0	0	0	0	0	0	0	0
47	0	0	0	0	0	0	0	0	0	0	0
48	0	0	0	0	0	0	0	0	0	0	0
49	0	0	0	0	0	0	0	0	0	0	0
50	0	0	0	0	0	0	0	0	0	0	0
51	0	0	0	0	0	0	0	0	0	0	0
52	0	0	0	0	0	0	0	0	0	0	0
53	0	0	0	0	0	0	0	0	0	0	0
54	0	0	0	0	0	0	0	0	0	0	0
55	0	0	0	0	0	0	0	0	0	0	0
56	0	0	0	0	0	0	0	0	0	0	0
57	0	0	0	0	0	0	0	0	0	0	0
58	0	0	0	0	0	0	0	0	0	0	0
59	0	0	0	0	0	0	0	0	0	0	0
60	0	0	0	0	0	0	0	0	0	0	0
61	0	0	0	0	0	0	0	0	0	0	0
62	0	0	0	0	0	0	0	0	0	0	0
63	0	0	0	0	0	0	0	0	0	0	0
64	0	0	0	0	0	0	0	0	0	0	0
65	0	0	0	0	0	0	0	0	0	0	0
66	0	0	0	0	0	0	0	0	0	0	0
67	0	0	0	0	0	0	0	0	0	0	0
68	0	0	0	0	0	0	0	0	0	0	0
69	0	0	0	0	0	0	0	0	0	0	0
70	0	0	0	0	0	0	0	0	0	0	0
71	0	0	0	0	0	0	0	0	0	0	0
72	0	0	0	0	0	0	0	0	0	0	0
73	0	0	0	0	0	0	0	0	0	0	0
74	0	0	0	0	0	0	0	0	0	0	0
75	0	0	0	0	0	0	0	0	0	0	0
76	0	0	0	0	0	0	0	0	0	0	0
77	0	0	0	0	0	0	0	0	0	0	0
78	0	0	0	0	0	0	0	0	0	0	0
79	0	0	0	0	0	0	0	0	0	0	0
80	0	0	0	0	0	0	0	0	0	0	0
81	0	0	0	0	0	0	0	0	0	0	0
82	0	0	0	0	0	0	0	0	0	0	0
83	0	0	0	0	0	0	0	0	0	0	0
84	0	0	0	0	0	0	0	0	0	0	0
85	0	0	0	0	0	0	0	0	0	0	0
86	0	0	0	0	0	0	0	0	0	0	0
87	0	0	0	0	0	0	0	0	0	0	0
88	0	0	0	0	0	0	0	0	0	0	0
89	0	0	0	0	0	0	0	0	0	0	0
90	0	0	0	0	0	0	0	0	0	0	0
91	0	0	0	0	0	0	0	0	0	0	0
92	0	0	0	0	0	0	0	0	0	0	0
93	0	0	0	0	0	0	0	0	0	0	0
94	0	0	0	0	0	0	0	0	0	0	0
95	0	0	0	0	0	0	0	0	0	0	0
96	0	0	0	0	0	0	0	0	0	0	0
97	0	0	0	0	0	0	0	0	0	0	0
98	0	0	0	0	0	0	0	0	0	0	0
99	0	0	0	0	0	0	0	0	0	0	0
100	0	0	0	0	0	0	0	0	0	0	0
101	0	0	0	0	0	0	0	0	0	0	0
102	0	0	0	0	0	0	0	0	0	0	0
103	0	0	0	0	0	0	0	0	0	0	0
104	0	0	0	0	0	0	0	0	0	0	0
105	0	0	0	0	0	0	0	0	0	0	0
106	0	0	0	0	0	0	0	0	0	0	0
107	0	0	0	0	0	0	0	0	0	0	0
108	0	0	0	0	0	0	0	0	0	0	0
109	0	0	0	0	0	0	0	0	0	0	0
110	0	0	0	0	0	0	0	0	0	0	0
111	0	0	0	0	0	0	0	0	0	0	0
112	0	0	0	0	0	0	0	0	0	0	0
113	0	0	0	0	0	0	0	0	0	0	0
114	0	0	0	0	0	0	0	0	0	0	0
115	0	0	0	0	0	0	0	0	0	0	0
116	0	0	0	0	0	0	0	0	0	0	0
117	0	0	0	0	0	0	0	0	0	0	0
118	0	0	0	0	0	0	0	0	0	0	0
119	0	0	0	0	0	0	0	0	0	0	0
120	0	0	0	0	0	0	0	0	0	0	0
121	0	0	0	0	0	0	0	0	0	0	0
122	0	0	0	0	0	0	0	0	0	0	0
123	0	0	0	0	0	0	0	0	0	0	0
124	0	0	0	0	0	0	0	0	0	0	0
125	0	0	0	0	0	0	0	0	0	0	0
126	0	0	0	0	0	0	0	0	0	0	0
127	0	0	0	0	0	0	0	0	0	0	0
128	0	0	0	0	0	0	0	0	0	0	0
129	0	0	0	0	0	0	0	0	0	0	0
130	0	0	0	0	0	0	0	0	0	0	0
131	0	0	0	0	0	0	0	0	0	0	0
132	0	0	0	0	0	0	0	0	0	0	0
133	0	0	0	0	0	0	0	0	0	0	0
134	0	0	0	0	0	0	0	0	0	0	0
135	0	0	0	0	0	0	0	0	0	0	0
136	0	0	0	0	0	0	0	0	0	0	0
137	0	0	0	0	0	0	0	0	0	0	0
138	0	0	0	0	0	0	0	0	0	0	0
139	0	0	0	0	0	0	0	0	0	0	0
140	0	0	0	0	0	0	0	0	0	0	0
141	0	0	0	0	0	0	0	0	0	0	0
142	0	0	0	0	0	0	0	0	0	0	0
143	0	0	0	0	0	0	0	0	0	0	0
144	0	0	0	0	0	0	0	0	0	0	0
145	0	0	0	0	0	0	0	0	0	0	0
146	0	0	0	0	0	0	0	0	0	0	0
147	0	0	0	0	0	0	0	0	0	0	0
148	0	0	0	0	0	0	0	0	0	0	0
149	0	0	0	0	0	0	0	0	0	0	0
150	0	0	0	0	0	0	0	0	0	0	0

$$r = +.737 \pm .027$$

COMPARISON TABLE FOR ACCURACY ON TERMAN FORM A WITH
 ACCURACY ON FORM B (HIGH GROUP I HIGH C)
 ACCURACY ON TERMAN FORM B (FIRST TRIAL)

ACCURACY ON TERMAN FORM A		25-29	30-34	35-39	40-44	45-49	50-54	55-59	60-64	65-69	70-74	75-79	80-84	85-89	90-94	+
	95	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	90	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	85	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	80	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	75	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	70	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	65	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	60	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	55	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	50	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	45	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	40	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	35	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
	30	•	•	•	•	•	•	•	•	•	•	•	•	•	•	•
f		2	2	10	11	10	10	10	12	11	11	7	9	11	11	195

$$\bar{r} = +.841 \pm .017$$

$$M_A = 0.633 \pm 0.007$$

$$M_B = 0.585 \pm 0.008$$

$$M_A - M_B = 0.048 \pm 0.011$$

$$\sigma_B = 0.1322 \pm 0.0054$$

$$\sigma_A = 0.1237 \pm 0.0051$$

$$\sigma_B - \sigma_A = 0.0085 \pm 0.007$$

∴ gain is significant

∴ diff. is not significant

APPENDIX B

Theorems Relating to Correlation

Notation:

 X, Y, Z -- variables from arbitrary origins x, y, z -- variables from respective means M_x --- means of the variables X , --- σ_x ---- standard deviations of X , ---- V_x ----- Pearson's coefficient of variability, $\frac{100\sigma}{M}$ r_{xy} ---- product moment correlation Σ ----- sum of such quantities as ----- N ----- frequency of the population a, b, c -- constants

Theorem 1. The correlation between two variables is the same as that between any two linear functions of each of them i.e.

$$(1) \quad r_{(aX+b)(cY+d)} = r_{xy} \quad (a, c \neq 0)$$

Transferring the variables to their respective means,

$$\begin{aligned} r_{(aX+b)(cY+d)} &= r_{(ax)(cy)} \\ &= \frac{\sum axcy}{\sqrt{\sum a^2x^2 \sum c^2y^2}} \\ &= \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}} \\ &= r_{xy} \end{aligned}$$

Theorem 2. The correlation between two ratios $\frac{X}{Y}$ and $\frac{Z}{W}$ is given by the formula,

$$(2) \quad r_{\frac{X}{Y} \frac{Z}{W}} = \frac{r_{xz} V_x V_z + r_{yw} V_y V_w - r_{xw} V_x V_w - r_{yz} V_y V_z}{\sqrt{[V_x^2 - 2r_{xy} V_x V_y + V_y^2]} [V_z^2 - 2r_{zw} V_z V_w + V_w^2]}$$

The means and standard deviations of $\frac{X}{Y}$ are given by

equations (9) and (10) in Yule,^a

$$(a) \quad M_{\frac{x}{Y}} = \frac{M_x}{M_Y} (1 - R_{xy} V_x V_Y + V_Y^2)$$

$$(b) \quad \sigma_{\frac{x}{Y}}^2 = \frac{M_x^2}{M_Y^2} (V_x^2 - 2R_{xy} V_x V_Y + V_Y^2)$$

The required correlation $R_{\frac{x}{Y} \frac{z}{W}}$ will then be given by

$$\begin{aligned} N R_{\frac{x}{Y} \frac{z}{W}} \sigma_{\frac{x}{Y}} \sigma_{\frac{z}{W}} &= \sum \left(\frac{x}{Y} - M_{\frac{x}{Y}} \right) \left(\frac{z}{W} - M_{\frac{z}{W}} \right) \\ &= \sum \left(\frac{x}{Y} \frac{z}{W} \right) - N M_{\frac{x}{Y}} M_{\frac{z}{W}} \\ &= \frac{M_x M_z}{M_Y M_W} \sum \left(1 + \frac{x}{M_x} \right) \left(1 + \frac{y}{M_Y} \right) \left(1 + \frac{z}{M_z} \right) \left(1 + \frac{w}{M_W} \right) - N M_{\frac{x}{Y}} M_{\frac{z}{W}} \end{aligned}$$

Expanding, neglecting terms higher than the second degree and substituting from (a) and (b) gives (3).

Formula (2) gives satisfactory results for fairly long series, but for very short ones, considerable error occurs due no doubt to neglecting the higher powers in the expansions of the binomials.

Theorem 3. If $X = aY + bZ$, the partial correlations between the variables may be written,

$$(3) \quad \begin{aligned} R_{xy.z} &= +1. \\ R_{xz.y} &= +1. \\ R_{yz.x} &= -1. \end{aligned}$$

as is evident by inspection.

Example: Since $A = R + W$, $R_{RW} = -1$, if $A = \text{const.}$

Theorem 4. The correlation between a variable and the sum of n others is given by the formula^a,

$$(4) \quad r_{(x_1+x_2+\dots+x_n)y} = \frac{r_{x_1y}\sigma_{x_1} + r_{x_2y}\sigma_{x_2} + \dots + r_{x_ny}\sigma_{x_n}}{\sqrt{\sigma_{x_1}^2 + \sigma_{x_2}^2 + \dots + \sigma_{x_n}^2 + 2(r_{x_1x_2}\sigma_{x_1}\sigma_{x_2} + \dots + \frac{n(n-1)}{2} \text{ terms})}}$$

(this expression is a special case of Spearman's General Formula)

$$r_{(x_1+x_2+\dots+x_n)y} = \frac{\sum x_1y + \sum x_2y + \dots + \sum x_ny}{N\sigma_{(x_1+x_2+\dots+x_n)}\sigma_y}$$

Since $\sigma_{(x_1+x_2+\dots+x_n)} = \sqrt{\sigma_{x_1}^2 + \sigma_{x_2}^2 + \dots + \sigma_{x_n}^2 + 2(r_{x_1x_2}\sigma_{x_1}\sigma_{x_2} + \dots + \frac{n(n-1)}{2} \text{ terms})}$ and $\sum x_1y = N r_{xy}\sigma_x\sigma_y$, the right hand member above reduces to the expression in (4)

Corollary. If the standard deviations σ_x are all equal

$$r_{(x_1+x_2+\dots+x_n)y} = \frac{r_{x_1y} + r_{x_2y} + \dots + r_{x_ny}}{\sqrt{n + 2(r_{x_1x_2} + r_{x_1x_3} + \dots + \frac{n(n-1)}{2} \text{ terms})}}$$

Theorem 5. The correlation between the sum of n variables and n other variables is given by the formula^b

$$(5) \quad r_{(x_1+x_2+\dots+x_n)(x'_1+x'_2+\dots+x'_n)} = \frac{r_{x_1x'_1}\sigma_{x_1}\sigma_{x'_1} + r_{x_1x'_2}\sigma_{x_1}\sigma_{x'_2} + \dots + n^2 \text{ terms}}{\sqrt{[\sigma_{x_1}^2 + \sigma_{x_n}^2 + 2(r_{x_1x_2}\sigma_{x_1}\sigma_{x_2} + \dots)]} [\sigma_{x'_1}^2 + \sigma_{x'_n}^2 + 2(\dots)]}$$

The proof is similar to (4)

Corollary 1. If the standard deviations σ_x are all equal

$$(6) \quad r_{nn} = \frac{r_{x_1x'_1} + r_{x_1x'_2} + \dots + n^2 \text{ terms}}{\sqrt{[n + 2(r_{x_1x_2} + \dots + \frac{n(n-1)}{2} \text{ terms})] [n + 2(r_{x'_1x'_2} + \dots + \frac{n(n-1)}{2} \text{ terms})]}}$$

where r_{nn} denotes the left hand member of (5).

a. C. Spearman, British Journal of Psychology, 1913, Vol.V, p.417.

b. C. Spearman, loc. cit.

Corollary 2. If the correlations ρ_{xx} are all equal,
equation (6) may be written

$$(7) \quad \rho_{nn} = \frac{n \rho_{xx}}{1 + (n-1) \rho_{xx}}$$

This is Brown's Theorem, but as shown above it is merely
 a special case of Spearman's General Formula.

10. 11. 1944

11. 11. 1944



77 20.55 1

BEWARE OF LOSS
REPLACEMENT PRICE
OVER FIFTEEN CENTS
PER PAGE
Receipt for return
given if requested.

LIBRARY OF CONGRESS



0 021 337 968 9